



# **DOCA Pipeline Language Services Guide**

# Table of contents

DPL Release Notes	13
DPL System Overview	17
DPL Runtime Service	25
Container Deployment	31
Service Configuration	36
DPL Development Container	43
DOCA Target Architecture	43
P4 Language Support in DPL	114
DPL Installation Guide	138
Compiling DPL Applications	141
Loading DPL Applications	145
Sample DPL Applications	147
Hello Packet Example	149
GTP Parsing Example	151
GTP Tunnel Encapsulation Example	155
Geneve TLV Parsing Example	163
VXLAN Tunnel Gateway Example	170
Connection Tracking Example	176
Host-based Networking Example	184

The DPL Services consist of 2 packages that form the DPL solution. The services are provided as containers and are deployed separately.

See the following sections on each service:

- [DPL System Overview](#) to get a high level understanding of the components that make up the DPL Services
- [DPL Runtime Service](#) to understand how to deploy and configure the backend DPL Runtime Service, that interacts with the hardware.
- [DPL Development Container](#) to learn about the DPL language, the compiler tools and methodology building DPL programs

See also:

- [DOCA Pipeline Language Developer Tool](#) to learn about the various tools and methodology for debugging DPL programs

### Info

For questions, comments, and feedback, please contact us at [DOCA-Feedback@exchange.nvidia.com](mailto:DOCA-Feedback@exchange.nvidia.com).

## DOCA Pipeline Language Model

This page outlines the DOCA Pipeline Language (DPL) approach to packet processing programmability for NVIDIA® BlueField®. DPL introduces a software development model based on a domain-specific programming language (DSL), supported by a set of DOCA services.

For in-depth details, refer to the [DOCA Pipeline Language Services Guide](#).

# Introduction

DPL is derived from the [P4-16 language specification](#).

P4 is an open-source, domain-specific programming language (DSL) designed for programming and customizing network data planes. It provides a high-level abstraction for programmable packet processing, allowing developers to add, modify, and extend networking functionalities.

For fixed-function devices, P4 serves as a documentation tool, offering a structured description of the data plane's functional blocks.

## Key Features of P4

- High-level abstraction – Simplifies programming of complex network data planes with clear and concise syntax
- Programmable packet processing – Enables customization of packet processing and traffic management
- Documentation of fixed functions – Offers a standardized method for documenting the fixed functional blocks of network devices

## P4 Compiler

The P4 compiler (p4c) is a critical component in the P4 ecosystem. It automatically generates the data plane program and a corresponding control plane interface, ensuring seamless coordination between the data plane and control plane.

Key benefits of the P4 compiler:

- Automatic generation – Streamlines development by automatically generating essential components and optimizing resource usage
- Custom pipeline behavior – Allows developers to extend data plane functionality with customized pipeline behaviors
- Dynamically loadable pipelines – Supports hot-swappable pipelines, enabling updates without rebuilding or redeploying an entire application
- Control plane integration – Facilitates communication between the data plane and control plane via an open-source API, ensuring effective management of customized

pipelines

## Focus on NVIDIA's DOCA Pipeline Language

The remainder of this document focuses on NVIDIA's implementation of the DOCA Pipeline Language (DPL). While DPL's syntax is derived from P4-16, its pipeline semantics align with NVIDIA's DPU pipeline architecture rather than standard P4 execution models.

For example, while P4 semantics imply a staged pipeline based on a feed-forward RMT (Reconfigurable Match-Action Table) architecture, NVIDIA's DPU architecture follows a run-to-completion dRMT (disaggregated RMT) model, offering greater flexibility and enhanced capabilities.

## DPL Highlights

The DPL introduces a unique programming paradigm distinct from traditional SDKs, APIs, libraries, drivers, or utilities. It is a specialized programming language with a runtime system, designed for rapid development, testing, and deployment of packet processing pipelines. DPL is provided as a ready-to-use, customizable solution under [DOCA Services](#).

### Key Features of DPL

- [DPL Services](#) – A system-level solution that includes a compiler, runtime agent, and debugging tools, enabling rapid programming of the DPU pipeline
- Optimized for NVIDIA devices – Specifically designed and fine-tuned for programming network data planes on NVIDIA hardware
- Advanced networking functionality – Leverages DPL's capabilities to enhance and extend networking features on NVIDIA DPUs
- Comprehensive documentation – Provides detailed descriptions of BlueField's fixed functional blocks within the DPU data plane

### Developer Resources

The DPL programming guide serves as a comprehensive resource for developers looking to harness DPL for programming network data planes. By utilizing the DPL `p4c` compiler

and the P4-16 specification, developers can:

- Enhance network device functionality and efficiency
- Meet the evolving demands of modern network infrastructures
- Ensure seamless integration and optimization within NVIDIA's DPU ecosystem

## Prerequisites

To effectively develop with DPL, readers should be familiar with the fundamentals of P4 and DPL. Language specifications, runtime APIs, and tutorials are available at [P4 GitHub Repository](#).

The DPL compiler can run on any Linux OS that supports Docker.

## Development Environment Requirements

To set up the development environment, the following components are required:

- Host Computer – Ubuntu 22.04 or later with Docker installed (required for the DOCA development container)
- Server with Root/Hypervisor Access – Required to install the DPL Runtime Service package
- One or More BlueField-3 Devices – Installed in the target server for DPL execution

## Suggested Workflow

The suggested workflow is as follows:

### 1. Coding

1. Using the DPL programming guide and sample applications, the developer creates a DPL program remotely.

2. The program is compiled using `dplp4c`, iterating until it successfully produces a binary.

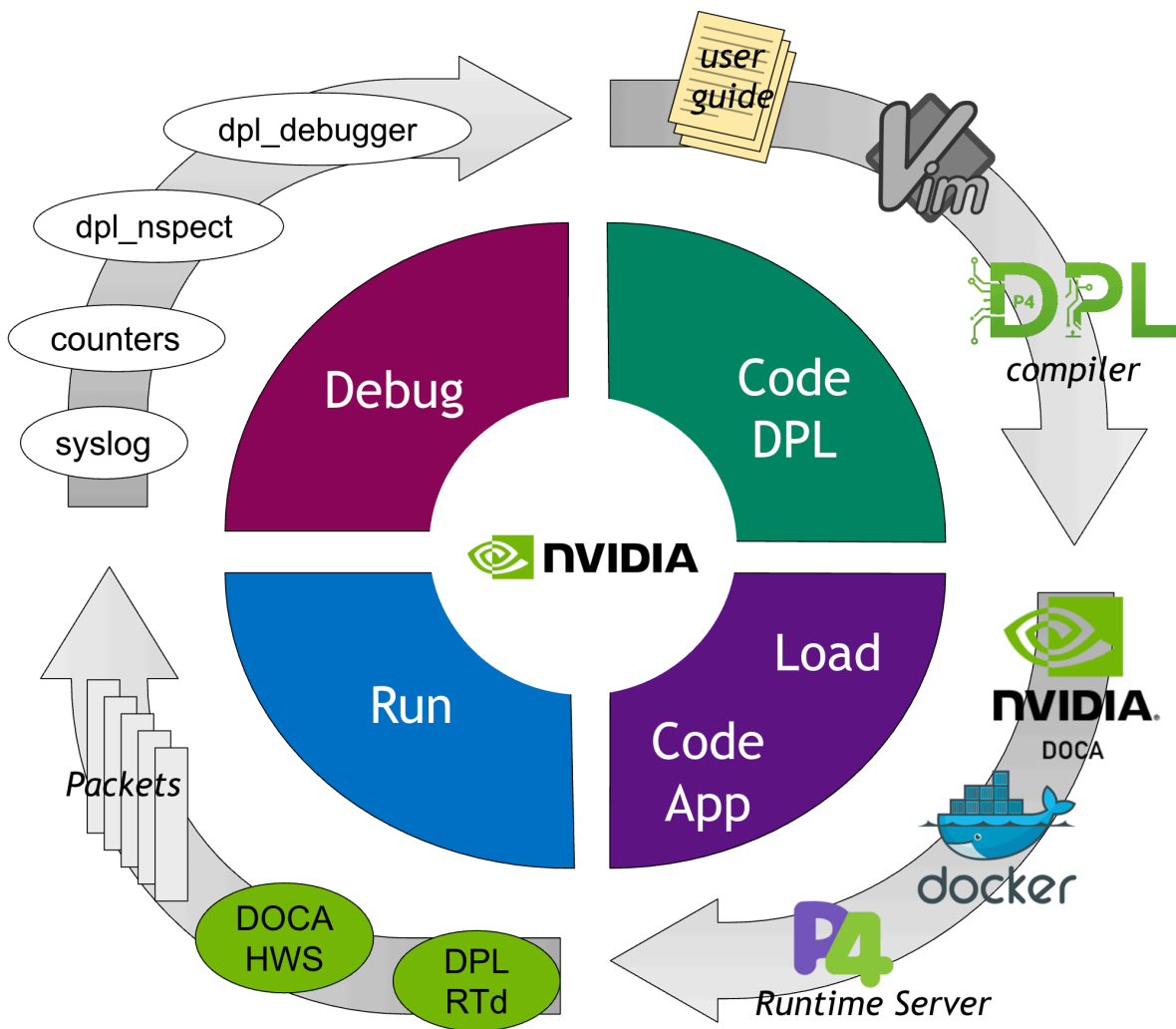
## 2. Loading

1. The compiled binary is transferred to the BlueField system.
2. Using the P4Runtime API (via an open-source or proprietary P4Runtime controller), the pipeline is sent from the remote machine to the DPL Service running on BlueField.
3. The user checks for P4Runtime error messages.

## 3. Running

1. The user inspects logs for any DPL Service error messages.
2. The `dpl_nspect` tool is used to verify that P4 tables and entries are present in the hardware.
3. The `dpl_pipeline_debugger` provides insights into the packet processing pipeline, showing the state of packets and their metadata.

This process is repeated until the DPL application is fully verified.



## DPL Programming Model

P4, and by extension DPL, is a domain-specific language (DSL) designed for programming network data planes. It enables customized packet processing, allowing developers to define how packets are handled at different pipeline stages.

However, P4 programs are not universally portable across different architectures. Instead, they are typically compatible within the same target architecture family.

The BlueField programmable pipeline follows a hybrid model that leverages both hardware and software processing capabilities. It consists of three main stages:

1. Parsing
2. Match-action processing (Steering)



### 3. Forwarding database (FDB)

## Parsing

The BlueField native parser is the first stage of the packet processing pipeline. It is responsible for identifying and extracting packet headers, progressing through the protocol stack until the entire frame is parsed.

Key features:

- Predefined protocol headers and standard transitions based on IETF specifications
- On-demand reparsing at any stage (eliminating the need for reinjection or a final deparser stage)

## Flex Parsing

Flex parsing allows developers to integrate custom protocol headers into BlueField's hardware parsing engine. It consists of four components:

- Flex Arc In – Defines the transition from a native header to a Flex header
- Flex Header – Specifies header characteristics such as length and next protocol location
- Flex Sampler – Extracts specific bytes from the hardware, enabling their use in control blocks or table keys
- Flex Arc Out – Defines the transition from a Flex header back to a native header (or another Flex header)

The DPL compiler automatically generates Flex parsing components based on the developer's defined parse nodes and transitions.

## Operational Mode

- The DPL parser operates in a hybrid mode with a default native parser
- The compiler automatically integrates native headers and fields into DPL constructs
- The Flex parse graph consists of:
  - Nodes (either native or Flex)

- Arcs (transitions)
- Samplers for custom parsing operations

This design eliminates the need to redefine and reimplement standard IETF protocols and headers.

## Match-Action Processing (Steering)

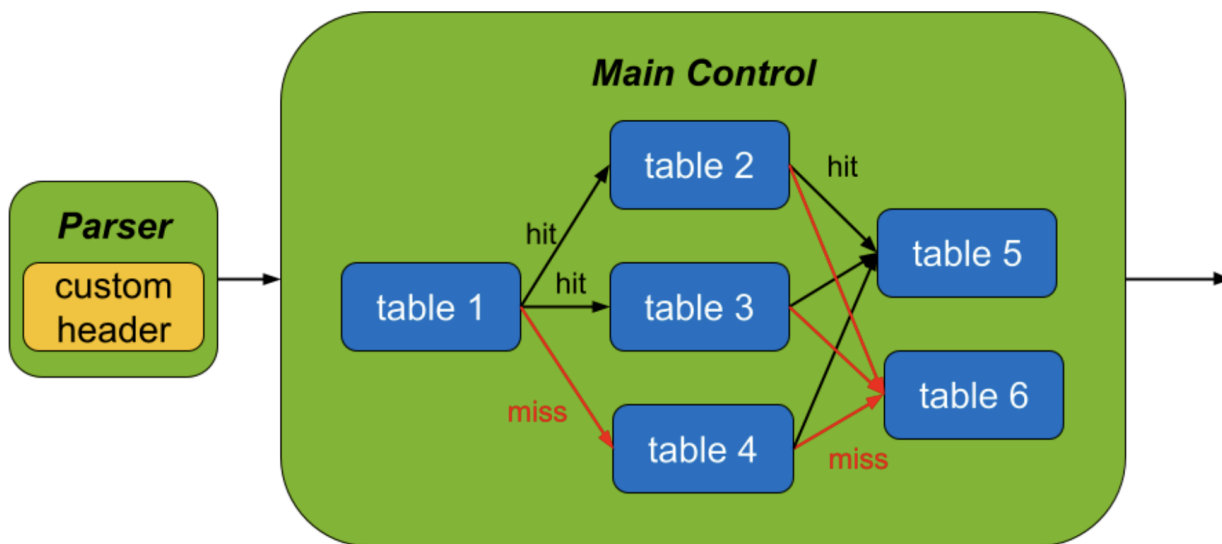
After parsing, packet processing decisions are made based on Match-Action tables, commonly referred to as "Steering".

Key Features:

- Match fields – Define packet attributes for classification (e.g., source/destination MAC, VLAN, IP, and protocol headers)
- Tables – Store rules for packet handling and decision-making
- Actions – Define processing rules (e.g., forwarding, header modification, dropping packets)
- Programmability – Allows dynamic updates to match-action rules based on network conditions
- Efficient processing – Packet handling occurs directly in hardware, reducing latency
- P4Runtime integration – DPL tables are populated via the P4Runtime API, supporting SDN controllers

### **Info**

In the documentation, flow tables may also be referred to as P4 tables.



## Forwarding Database

The Forwarding Database (FDB) is the final stage within the embedded switch (eSwitch). It is responsible for:

- Storing and managing MAC addresses
- Ensuring efficient packet forwarding based on network topology
- Maintaining records of port locations for destination-based forwarding

The FDB enables accurate and efficient packet routing within the network infrastructure.

## BlueField DPU Pipeline Behavior

The BlueField pipeline is designed for flexibility, allowing developers to customize packet processing to meet specific application needs.

Key characteristics:

- Extended parser support – Developers can expand the native parser using Flex Parsing
- Immediate execution model – No deferred actions; all modifications take effect immediately

- Mid-pipeline reparsing – Packet headers are reparsed immediately after modification, ensuring correct metadata updates
- No deparser control in TA – Unlike traditional architectures, BlueField does not require a separate deparser step

### Info

For example, after an encapsulation action, the changes are immediately visible to the next processing table.

## DPL Services

Rather than providing a traditional SDK or driver-level APIs, DPL offers a high-level services-based approach to programming the DPU pipeline.


The DPL Services consist of 2 packages that form the DPL solution. The services are provided as containers and are deployed separately.

See the following sections on each service:

- [DPL System Overview](#) to get a high level understanding of the components that make up the DPL Services
- [DPL Runtime Service](#) to understand how to deploy and configure the backend DPL Runtime Service, that interacts with the hardware.
- [DPL Development Container](#) to learn about the DPL language, the compiler tools and methodology building DPL programs

See the additional sections:

- [DOCA Pipeline Language Model](#) provides a brief overview of the DPL developer workflow and the DPL programming model.
- [DOCA Pipeline Language Developer Tool](#) to learn about the various tools and methodology for debugging DPL programs

 **Info**

For questions, comments, and feedback, please contact us at [DOCA-Feedback@exchange.nvidia.com](mailto:DOCA-Feedback@exchange.nvidia.com).

---

# DPL Release Notes

## Changes and New Features

- First release (GA)

## Capabilities and Limitations

### P4 Runtime

- The size of any of P4 table and counter object must be a power-of-2 (smaller than UINT32\_MAX)
- A P4 table of size N allows for "N-1" regular entries and 1 default entry
- P4 Controller support:
  - Only one P4 Controller can be connected to the DPL Runtime daemon at a given time.
- RPC messages support:
  - Supported RPC messages:
    - Write RPC:
      - Only CONTINUE\_ON\_ERROR atomicity is supported
      - Batching is supported
      - Supported entities:
        - TableEntry
        - CounterEntry
        - DirectCounterEntry
        - MeterEntry

- DirectMeterEntry
- INSERT operation:
  - Supports regular entries only
- DELETE operation:
  - Supports regular entries only
- MODIFY operation:
  - Supports Default entry only
  - Supports indirect and direct counter entry
  - Supports indirect and direct meter entry

Mode	max cir/pir	max cburst/pburst
<b>BYTES</b>	2550000000000	80000000
PACKETS	1992187500	625000

\* In packets mode packet = 128 bytes

- Read RPC:
  - Batching is supported
  - Supported entities:
    - TableEntry
    - CounterEntry
    - DirectCounterEntry
    - MeterEntry
    - DirectMeterEntry
- SetForwardingPipelineConfig RPC
- GetForwardingPipelineConfig RPC

- Unsupported RPC messages:
  - Capabilities RPC

## DPL Nspect

- `dpl_nspect graph` command – default graph type (`pipeline`) is not supported. Therefore, the `--type` argument is mandatory.
- `dpl_nspect graph --help` – shows the available graph types.

## Bug Fixes

N/A

## Known Issues

The following are known limitations of the DPL Runtime daemon service.

Ref #	Issue
4294992	Description: If the configuration file refers to interfaces that don't exist, the error is only reported when loading a program and is not friendly to the user
	Workaround: fix the configuration files or create the missing VFs or SFs
	Keywords: dpl_rtd
	Reported in version: N/A
4320688	Description: The dpl_rtd service may crash if adding constant table entries defined in the program fails.
	Workaround: No action needed, the dpl_rtd service will automatically restart.
	Keywords: dpl_rtd



Ref #	Issue
	Reported in version: N/A
1141	Description: Debugging packets from the second wire port P1 is currently not supported
	Workaround: Use wire port P0 for debugging packets.
	Keywords: debug
	Reported in version: N/A

---

# DPL System Overview

This section outlines the structure and capabilities of the NVIDIA® BlueField® programmable pipeline, focusing on its implementation within the DOCA Pipeline Language (DPL) framework. The syntax from the P4-16 language enables users to describe the behavior of the packet processing pipeline, while DPL provides an NVIDIA platform-specific target architecture (TA) optimized for the BlueField networking platform's hardware (DPU or SuperNIC).

## Introduction

The [P4 Language](#) is a powerful and flexible programming language designed to define advanced network behavior. DOCA Programming Language (DPL) is based on P4-16, reusing its syntax to provide a familiar programming interface. However, some of the semantics, programmability features and behavioral models are DOCA specific. Unlike traditional approaches that rely on APIs, DPL enables developers to specify packet-handling rules directly through code. These rules determine how packets are processed, including header parsing, matching logic, and actions such as forwarding, dropping, or modifying packets. DPL offers the flexibility to work with existing header formats or define custom headers. Developers can then use these header fields for matching and manipulation, providing extensive control over packet processing. Compared to implementing equivalent pipelines using C-based APIs, developing DPL programs is faster and more efficient. Modifications are simpler to make, and the DPL compiler automatically optimizes the pipeline, allowing developers to focus on correctness rather than performance.

When a DPL program is executed, runtime control is achieved using [P4Runtime](#), an open-source protocol that allows dynamic updates to rules without recompiling the software stack. This enables real-time reconfiguration of packet processing pipelines. While any P4Runtime compliant SDN controller can be used with DPL, some vendor specific extensions will be included in the future.

The **DPL Solution** provides a comprehensive toolkit for compiling, executing, monitoring, and debugging DPL programs. It includes a library of commonly used header types and consists of three parts: a [compiler](#), a [runtime service](#) and the [developer tools](#). The service runs on NVIDIA BlueField DPUs, while the development tools can operate on any system, such as a laptop. Together, these components facilitate efficient development and deployment of DPL-based applications on BlueField platforms.

# BlueField Network Configuration

BlueField-3 is a powerful cloud infrastructure processor that revolutionizes data center operations by offloading, accelerating, and isolating software-defined networking, storage, security, and management functions. In the context of processing packets in the data path, BlueField-3 leverages the NVIDIA® ConnectX® ASIC for high-speed networking capabilities, enabling efficient packet processing and data transfer. The Arm subsystem within the BlueField-3 provides programmable processing power for network functions and management tasks, enhancing overall system flexibility and performance. Additionally, BlueField-3 works in tandem with the host CPU to offload networking and security tasks, freeing up CPU resources to focus on running applications, thereby optimizing data center performance and efficiency.

BlueField-3 supports two DOCA device types:

- Local device – this is an actual device exposed in the local system (BlueField or host) and can perform DOCA library processing jobs. This can be a PCIe physical function (PF) virtual function (VF) or scalable function (SF)
- Netdev representor – this is a representation of a local device. The represented local device is usually on the host (except for SFs) and the representor is always on the BlueField side (a proxy on the BlueField for the host-side device). The representors map each one of the host side physical and virtual functions where it:
  - Serves as the tunnel to pass traffic for the virtual switch or application running on the Arm cores to the relevant PF or VF on the Arm side.
  - Serves as the channel to configure the embedded switch with rules to the corresponding represented function.

These representors are used as the virtual ports that connect to OVS and any other virtual switch running on the Arm cores.

When in embedded CPU function (ECPF) ownership mode (also called DPU mode), there are 2 representors for each of the BlueField's network ports: one for the uplink, and another one for the host side PF (the PF representor created even if the PF is not probed on the host side). For each one of the VFs created on the host side a corresponding representor is created on the Arm side. The naming convention for the representors is as follows:

- Uplink representors – `p<port_number>`
- PF representors – `pf<port_number>hpf`

- VF representors – `pf<port_number>vf<function_number>`
- SF representors - `en3f0pf0sf1`  
`en<pcie_domain>f<pcie_function>pf<port_number>sf<function_number>`

Traffic from the x86 host server hosting the BF3 to an external host will go via BlueField's Arm running OVS-DOCA. See the section BlueField DPU Mode under [DOCA Switching](#) for further details.

### **Note**

The MTU of host functions (PF/VF) must be smaller than the MTUs of both the uplink and corresponding PF/VF representor. For example, if the host PF MTU is set to 9000, both uplink and PF representor must be set to above 9000.

## DPL Architecture

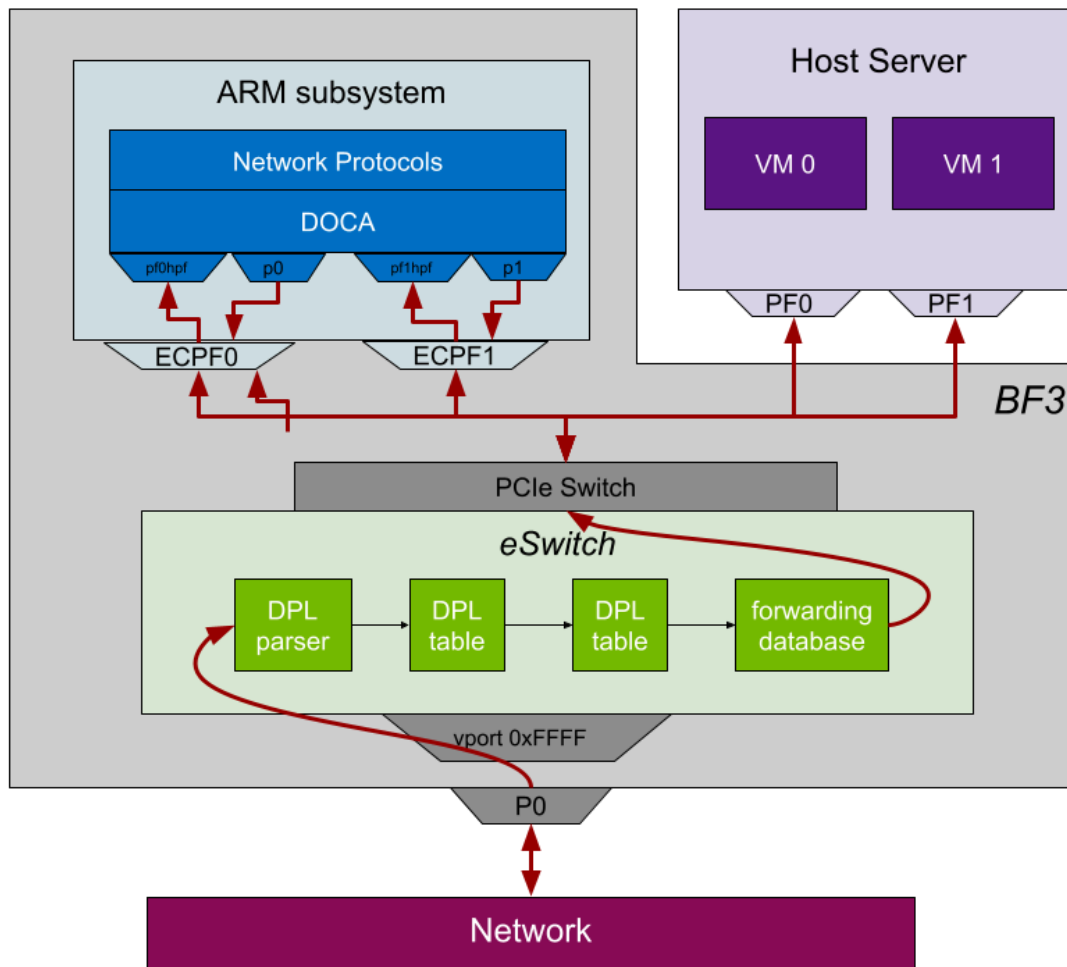
In DPU mode, the NIC resources and functionality are owned and controlled by the embedded Arm subsystem. All network communication to the host flows through a virtual switch control plane hosted on the Arm cores, and only then proceeds to the host. While working in this mode, BlueField is the trusted function managed by the data center and host administrator—to load network drivers, reset an interface, bring an interface up and down, update the firmware, and change the mode of operation on the BlueField DPU.

A network function is still exposed to the host, but it has limited privileges. In particular:

1. The driver on the host side can only be loaded after the driver on BlueField has loaded and completed NIC configuration.
2. All ICM (Interface Configuration Memory) is allocated by the ECPF and resides in BlueField's memory.
3. The ECPF controls and configures the NIC embedded switch which means that traffic to and from the host (BlueField) interface always lands on the Arm side.

A deployed DPL program executes on the eSwitch component of the BlueField hardware. Basic DPL processing blocks consist of:

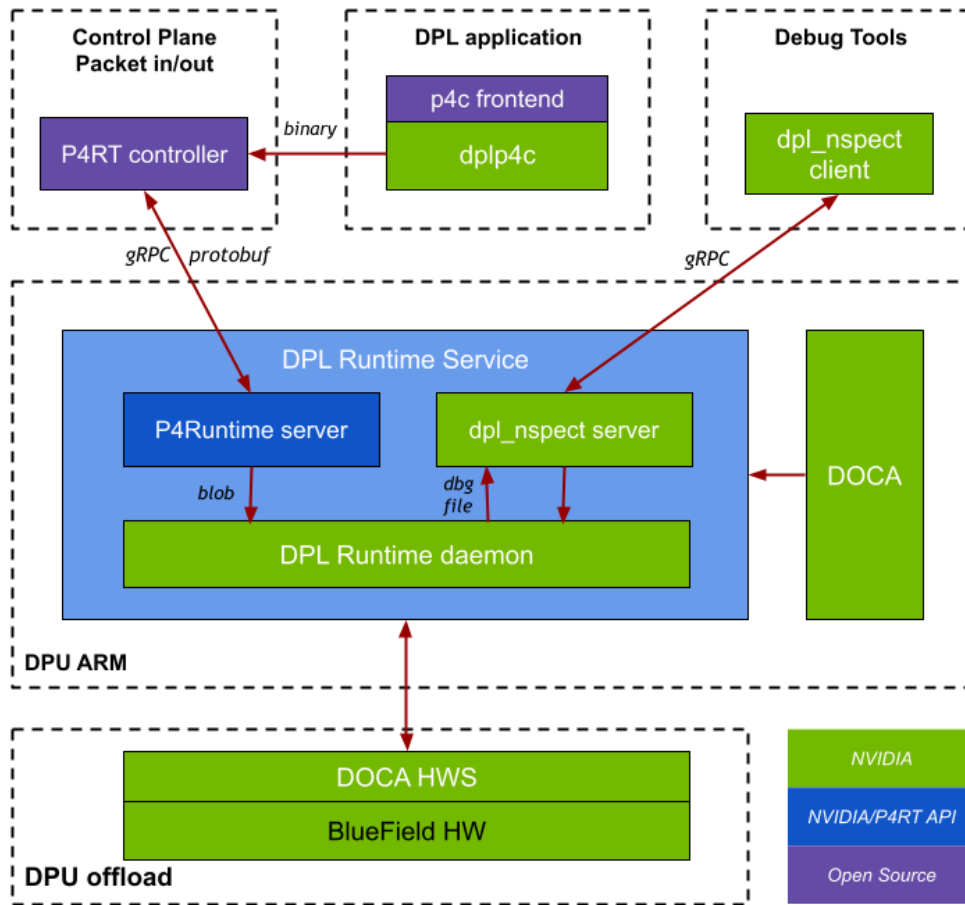
- A programmable packet parser based on a native parse graph
- A pipeline of match/action flow tables
- Actions on packets based on the outcome of a lookup. Actions include packet modification, encapsulation, payload encryption, forwarding actions and packet duplication.



- `pf0hpf` and `pf1hpf` are the representors facing the x86 host
- `pf0` and `pf1` are the representors facing the network

## DPL Components

The DPL solution on the BlueField DPU consists of multiple components working in tandem. These components are outlined below and shown in the following diagram:



## DPL Compiler

The NVIDIA DPL Compiler (`dplp4c`) is the central component of the DPL ecosystem. The compiler is typically used from a remote host CPU equipped with Ubuntu 22.04, rather than locally on BlueField. Its core functions include processing and analyzing DPL programs, with the objective of crafting an optimized pipeline compatible with BlueField hardware capabilities. This process culminates in the generation of two critical output files:

- P4info file – This protobuf text file is designed to be target-agnostic, facilitating the integration with P4Runtime controllers by supporting P4Runtime APIs (version 1.3).
- Device configuration file – This binary `*.dplconfig` file is tailored specifically to the BlueField target device. It contains essential data structures that bridge the gap between P4Runtime (P4RT) commands and DOCA HWS driver objects, ensuring seamless translation and execution of DPL programs on BlueField hardware.

An optional debug information file is produced that is used by the DPL debug tools. For detailed specifications on the P4 language, please refer to the [P4-16 v 1.2.4 documentation](#).

## **P4Runtime Controller**

The P4Runtime Controller is a dynamic management tool designed for networks utilizing P4-programmable devices, enabling network operators to configure, control, and monitor these devices in real-time. This user-owned component must comply with the [P4Runtime Specification 1.3.0](#), ensuring standardized communication and functionality across different network environments. The controller provides the ability to dynamically update packet processing rules, deploy custom network functions, and adapt to changing network conditions.

## **DPL Nspect Client**

DPL Nspect is a dedicated debugging tool tailored for diagnosing and troubleshooting DPL applications on BlueField devices. This standalone utility assists developers in understanding how their DPL program's concepts are implemented on actual hardware. It simplifies various debugging tasks, such as examining table contents and analyzing debug packets to trace issues back to their source. DPL Nspect can be installed on any host machine that has network access to the BlueField's management IP interface, making it a versatile and essential tool for developers working with DPL-programmed BlueField devices.

## **DPL Runtime Service**

The DPL Runtime service process fulfills various roles within BlueField's system:

- Manages the gRPC server and protobuf messages from the P4Runtime controller
- Processes compiler-generated P4info files and binary artifacts
- Responds to controller requests and provides BlueField status updates
- Handles P4 packet-in and packet-out data streams

- Bridges platform-independent P4Runtime APIs with platform-dependent functions using DOCA HWS Driver APIs
- Loads compiler binary artifacts to realize the DPL pipeline in BlueField hardware
- Maintains a database resource mapping between P4Runtime and DOCA HWS Driver APIs along with DPL parse graph information
- Manages the gRPC server for the DPL developer tools

## **P4Runtime Server**

The P4Runtime Server enables communication between the controller(s) and the data plane, allowing for dynamic updates to packet processing rules, querying device configurations, and managing P4 entities declared in the P4Info metadata. The P4 Runtime Server is implemented as a gRPC server, which binds an implementation of auto-generated client and server stubs based on the `p4runtime.proto` Protobuf file. It listens on TCP port 9559 by default and allows for program-independent control of P4 targets, ensuring a target-independent and protocol-independent approach to runtime control. The software component responsible for implementing the server functionality on BlueField is the DPL Runtime service. The server is the northbound component that processes the P4 Runtime messages by deserializing the messages and dispatching them to the `dpl_rtd` for execution.

## **DPL Nspect Server**

The DPL Nspect server is an internal component of the DPL Runtime Service, and it provides services to the DPL Nspect and DPL Debugger tools. The server implements a dedicated gRPC/protobuf-based server for the client to receive monitoring and debug information about the running DPL program.

## **DOCA HWS Driver**

The DOCA HWS Driver API provides the essential building blocks for creating the packet processing pipeline in hardware. It provides APIs for constructing flow tables with match criteria, actions, and flow control logic. The compiler utilizes these APIs to build user-defined pipelines efficiently. The DOCA HWS driver serves as the abstraction layer for programmable frameworks such as the DPL compiler and the underlying firmware and BlueField hardware.



## **eSwitch Management**

As a programmability solution, it is the developer's choice as to how the eSwitch administrator manages bridges, virtual ports, interfaces and global pipeline configuration. While Open vSwitch (OVS), a popular open-source virtual switch implementation (e.g. DOCA OVS) provides a flexible, multi-layer virtual networking switch, the user can configure interfaces via DPDK or DOCA as required by the application.

---

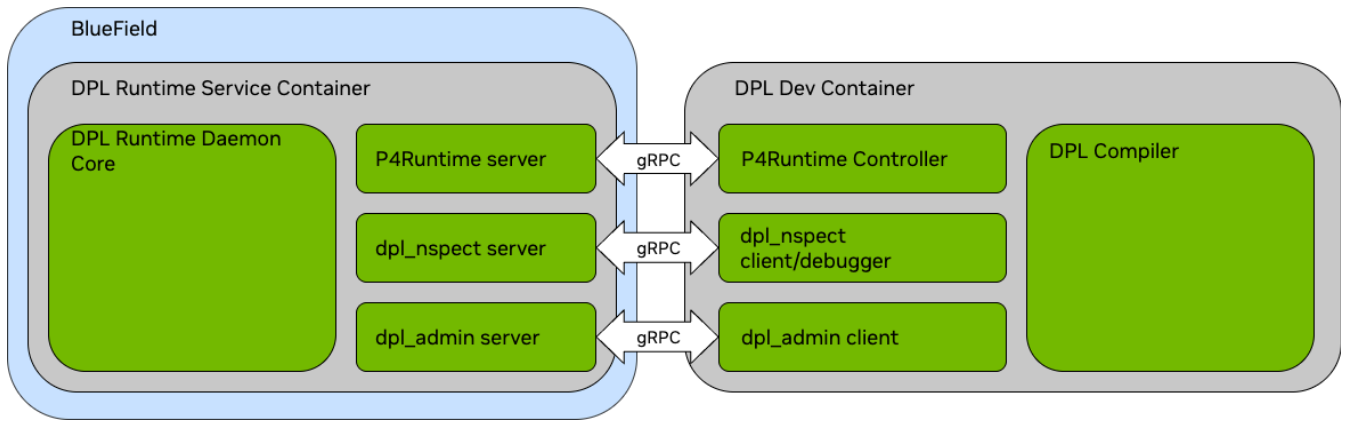
# DPL Runtime Service

## DPL Runtime Service

The DPL Runtime Service is an NVIDIA® BlueField® service that implements the backend functionality to manage and program the DPU datapath. It is broken down into 4 components:

- DPL Runtime Core. This component serves as the manager of the runtime system, accepting requests from the various servers and is responsible for managing resources and programming the hardware.
- dpl\_nspect server. This gRPC based server receives requests from the dpl\_nspect client application for debug information, and transmits debug packets to the DPL debugger.
- dpl\_admin server. This gRPC based server receives administrative requests about the daemon's core state and can control its configuration dynamically.
- P4Runtime server. This gRPC based server binds the P4Runtime protobuf interface to the underlying hardware driver APIs. The P4Runtime server listens on TCP port 9559, which is the port that has been allocated by IANA for the P4Runtime service. The server allows a P4 Controller to connect over gRPC so that it can set the `ForwardingPipelineConfig`, which installs and loads into hardware the compiled DPL program output and the associated P4Info metadata. Furthermore, the controller can query the target for the `ForwardingPipelineConfig` to retrieve the device config and the P4Info, as well as performing P4 table maintenance (which were defined in the DPL program source code).

## High level system illustration



## Requirements

### Supported Platforms

The following NVIDIA® BlueField® DPUs and SuperNICs are supported with DPL:

NVIDIA SKU	Legacy OPN	PSID	Description
900-9D3B6-00CV-AA0	N/A	MT_000000884	BlueField-3 B3220 P-Series FHHL DPU; 200GbE (default mode) / NDR200 IB; Dual-port QSFP112; PCIe Gen5.0 x16 with x16 PCIe extension option; 16 Arm cores; 32GB on-board DDR; integrated BMC; Crypto Enabled
900-9D3B6-00SV-AA0	N/A	MT_000000965	BlueField-3 B3220 P-Series FHHL DPU; 200GbE (default mode) / NDR200 IB; Dual-port QSFP112; PCIe Gen5.0 x16 with x16 PCIe extension option; 16 Arm cores; 32GB on-board DDR; integrated BMC; Crypto Disabled
900-9D3B6-00CC-AA0	N/A	MT_000001024	BlueField-3 B3210 P-Series FHHL DPU; 100GbE (default mode) / HDR100 IB; Dual-port QSFP112; PCIe Gen5.0 x16 with x16 PCIe extension option; 16 Arm cores; 32GB on-board DDR; integrated BMC; Crypto Enabled
900-9D3B6-00SC-AA0	N/A	MT_000001025	BlueField-3 B3210 P-Series FHHL DPU; 100GbE (default mode) / HDR100 IB; Dual-port QSFP112; PCIe Gen5.0 x16 with x16 PCIe extension option; 16 Arm cores; 32GB on-board DDR; integrated BMC; Crypto Disabled

The following NVIDIA® DPUs are supported with DOCA on the host:

- [Bluefield-3 devices](#)

## Hardware Prerequisites

For the the system requirements, see [DPU's hardware user guide](#).

## Software Prerequisites

This quick start guide assumes that an NVIDIA® BlueField® networking platform has been installed on a server. The minimum version requirements are described in the table below:

OS / Product	Version
ubuntu	22.04
bf-bundle	2.10.0-xx

## Installation Procedures

### Software Versions

The firmware version must be 32.44.1000 or higher. The firmware is included in the BFB bundle, which can be downloaded from the [NVIDIA Dev Zone](#). Select BlueField/BF-FW-Bundle/BFB. Installation example:

```
bfb-install --bfb bf-fwbundle-2.10.0-28_25.0-ubuntu-22.04-prod.bfb -  
-rshim rshim0
```

At a minimum [DOCA-networking](#) 2.10.0 should be installed on the host. For complete DOCA installation instructions, please see [DOCA Installation Guide for Linux](#).

The deployment guide has more information about this here: [Setup DPU Management access and update BlueField-Bundle](#)

## Firmware settings

The following firmware settings are needed

- FLEX\_PARSER\_PROFILE\_ENABLE=4
- PROG\_PARSE\_GRAPH=true
- SRIOV\_EN=1

(SRIOV\_EN=1 is only needed if you intend to use VFs)

For details on how to query and adjust these settings, please see [Using mlxconfig](#).

Some changes require a firmware reset, please see [mlxfwreset](#).

## How to deploy and run the DPL Runtime Service container

The DPL Runtime Service runs on the DPU of BlueField, which requires that BlueField is in DPU mode. for details, see this page: [BlueField Modes of Operation](#)

The DPL Runtime Service is deployed from [NGC](#), for detailed instructions, please refer to this page: [Container Deployment](#)

Both SR-IOV Virtual Functions (VFs) and Scalable Function (SFs) are supported. **VFs must be created on the host server before setting up the DPL Runtime Service in the DPU.**

Fetching the configuration files from NGC will create a directory named `dpl_rt_service_<version>`. Inside, you'll find `scripts/dpl_dpu_setup.sh`.

Running this script on the DPU will allow the usage of SR-IOV Virtual-Function interfaces and create the directory structure of the configuration files. In addition, the script will call the necessary `mlxconfig` commands needed for the DPL Runtime Service .

## DPL Runtime Service Configuration files

When the image is installed, containers can be created to run the DPL Runtime Service . Before that, however, you must create the configuration files that specify which network interfaces are usable by the DPL program and their IDs, along with other system parameters.

These configuration files are best created on the DPU file system to make them persistent and shared with the container file system. Containers by default have their own

separate file system unless specific paths are mounted in the container upon creation. Mounting the configuration path in the container is the preferred way to make sure the configuration is accessible to the DPL Runtime Service . For more details, see this page: [Service Configuration](#)

## Supported Port Forwarding

The following table lists out the possible ingress and egress ports for a given packet that is processed by a BlueField pipeline (DPU mode):

Ingress Port	Egress Port			
	Wire Port P0	PF0hpf	pf0vf_n	pf0vf_m
Wire port P0	Allowed	Allowed	Allowed	Allowed
pf0hpf	Allowed	Disabled	Allowed	Allowed
pf0vf_n	Allowed	Allowed	Disabled	Allowed
pf0vf_m	Allowed	Allowed	Allowed	Disabled

### Info

Anything that is allowed with SR-IOV Virtual Functions (VFs) in the table above is also allowed with Scalable Functions (SFs) .

To set the system in multiport e-switch mode, use the following command:

```
mlxconfig -d <pci> s LAG_RESOURCE_ALLOCATION=1
mlxfwreset -d <pci> -y -l 3 --sync 1 r
devlink dev param set <pci> name esw_multiport value 1 cmode
runtime
```

This mode will allow DPL to process packets on both wire ports P0 and P1. Note that the devlink command is not persistent across reboots.

# DPL Dev Container

Please refer to the [DPL Installation Guide](#) documentation.

## Implement your own P4Runtime Controller

The most powerful use-cases involve a controller that responds to live traffic and adapt the rules accordingly. This can be achieved by implementing your own P4Runtime client application that connects to the gRPC server of the running daemon. The controller provided in the DPL Dev Container is enough to load a DPL program but you can expect the best performance if you create an optimized client that is running directly on the Arm subsystem inside the BlueField DPU.

The necessary .proto files can be found here: [p4runtime/proto/p4 at v1.3.0 · p4lang/p4runtime · GitHub](#) In addition, it's possible to extract the files from a running DPL Runtime Service container. See inside the container in the directory `/opt/mellanox/third_party/dpl_rt_service/p4runtime/`.

This user guide does not explain how to compile and build the .proto files into C++ code and link them into an executable. Please see the information online for these open-source projects. For example, additional information about that can be found here: [C++ | gRPC](#)

## Troubleshooting

Kubelet logs can be viewed with command:

`sudo journalctl -u kubelet --since -5m` Installed (pulled) images can be observed with command: `sudo crictl images` Created pods can be observed with command: `sudo crictl pods`

Running containers pods can be observed with command: `sudo crictl ps` The log of the DPL Runtime Service is available in `/var/log/doca/dpl_rt_service/dpl_rtd.log`

Alternatively, the DPL Runtime Service logs may be observed from your development environment by using the [DPL Admin Control](#) tool.

Things to check:

- Did you create VFs in the host before setting up the DPU?

- Ensure that the BlueField-3 DPU is in `DPU MODE` mode. Refer to [BlueField Modes of Operation](#)
- Did you create configuration files? In the correct place?
- Did you follow the naming convention for the configuration files?
- Does the device ID in the configuration file match the file name?
- Did you get all the interface names correct?
- Is your firmware version up to date? (See [mlxup-mft](#))
- Make sure the link type is set to ETH in step 5 of the "Installing Software on Host" section in the [DOCA Installation Guide for Linux](#).

# Container Deployment

## Preparing the BlueField DPU

### Set BlueField to DPU Mode

BlueField must run in DPU mode to use the DPL Runtime Service . For details how to change modes, see here: [BlueField Modes of Operation](#).

### Determine Your BlueField Variant

Your BlueField may be Installed in a host server or it may be a standalone server.

If your BlueField is a standalone server, please ignore the parts that mention the host server or SR-IOV.

You may still use Scalable Functions (SFs) if your BlueField is a standalone server.

### Setup DPU Management Access and Update BlueField-Bundle



These pages provide detailed information about DPU management access and software installation and updates:

- [Host-side Interface Configuration - NVIDIA Docs](#)
- [BF-Bundle Installation and Upgrade](#)
- [NVIDIA DOCA Downloads | NVIDIA Developer](#)

Systems with a Host Server typically use RShim (i.e. the `tmfifo_net0` interface).

Standalone systems will have to use the OOB interface option for management access.

## Port Configuration

### Creating SR-IOV Virtual Functions (Host Server)

The first step to use SR-IOV is to create Virtual Functions (VFs) on the host server.

VFs can be created using the following sequence:

```
sudo -S # enter sudo shell
echo 4 > /sys/class/net/eth2/device/sriov_numvfs
exit # exit sudo shell
```

#### Info

Entering sudo shell rather than just issuing a single `sudo` command is necessary because otherwise the `sudo` applies only to the echo command and not the hosting shell and the redirection fails with "Permission denied"

This example creates 4 VFs under Physical Function eth2. Please adjust according to your needs.

If a PF already has VFs and you'd like to change the number of VFs, please set it to 0 before applying the new value.

## Scalable Functions (DPU)

For more information, see this: [BlueField Scalable Function User Guide](#)

If you create SFs, refer to their representors in the configuration file.

## Install the DPL Runtime Service on the DPU

## Pulling the Container Resources and Scripts from NGC

Start by downloading and installing the [ngc-cli](#) tools.

Fetch the configuration files from NGC, this will create a directory named `dpl_rt_service_<version>`.

e.g. `dpl_rt_service_v1.0.0-doca2.10.0`

Commands:

```
wget --content-disposition
https://api.ngc.nvidia.com/v2/resources/nvidia/ngc-
apps/ngc_cli/versions/3.58.0/files/ngccli_arm64.zip -O
ngccli_arm64.zip
unzip ngccli_arm64.zip
./ngc-cli/ngc registry resource download-version
"nvidia/doca/dpl_rt_service"
cd dpl_rt_service_v1.0.0-doca2.10.0
```

## Running the Preparation Script

Inside the directory with the scripts and YAML files that you pulled with the ngc-cli tool, you'll find `scripts/dpl_dpu_setup.sh`.

Running this script on the DPU (requires sudo) will allow the usage of SR-IOV Virtual-Function interfaces and create the directory structure of the configuration files in directory `/etc/dpl_rt_service`. In addition, the script will set "hugepages" and call the necessary `mlxconfig` commands to use DPL Runtime Service.

Run the following sequence of commands from the working directory you pulled with the ngc-cli tool:

```
chmod +x ./scripts/dpl_dpu_setup.sh
sudo ./scripts/dpl_dpu_setup.sh
sudo systemctl restart kubelet.service
sudo systemctl restart containerd.service
```

Restarting the services is necessary for the "hugepages" change to apply to them.

### Info

The following firmware settings are set by the setup script:

- FLEX\_PARSER\_PROFILE\_ENABLE=4
- PROG\_PARSE\_GRAPH=true
- SRIOV\_EN=1

## Edit the Configuration Files

Modify your configuration files as they are described here: [Service Configuration](#)

Important: you must create at least one device configuration under `/etc/dpl_rt_service/devices.d/`. It's advisable to start by making a copy of file `/etc/dpl_rt_service/devices.d/NAME.conf.template`.

e.g.

```
cp /etc/dpl_rt_service/devices.d/NAME.conf.template
   /etc/dpl_rt_service/devices.d/1000.conf
```

## Setting up the kubelet Pod

Now that everything is ready, copy the file `configs/dpl_rt_service.yaml` from the directory that you pulled with the `ngc-cli` into directory `/etc/kubelet.d`.

Please allow a few minutes for the image to be pulled and the pod to be started. you may check the progress with command `sudo journalctl -u kubelet --since -5m`, make sure to scroll down to see the latest log lines.

When the image is pulled, you will see it by using the command `sudo crictl images`.

When the pod is loaded, you will see it by using the command `sudo crictl pods`.

When the DPL Runtime Service is successfully running inside the pod, you will be able to find the log file in `/var/log/doca/dpl_rt_service/dpl_rtd.log`

## Recap, Full Command Sequence

```
wget --content-disposition
https://api.ngc.nvidia.com/v2/resources/nvidia/ngc-
apps/ngc_cli/versions/3.58.0/files/ngccli_arm64.zip -O
ngccli_arm64.zip
unzip ngccli_arm64.zip
```

```
./ngc-cli/ngc registry resource download-version
"nvidia/doca/dpl_rt_service"
cd dpl_rt_service_v1.0.0v1
chmod +x ./scripts/dpl_dpu_setup.sh
sudo ./scripts/dpl_dpu_setup.sh
sudo systemctl restart kubelet.service
sudo systemctl restart containerd.service

sudo cp /etc/dpl_rt_service/devices.d/NAME.conf.template
/etc/dpl_rt_service/devices.d/1000.conf
## Modify the configuration file /etc/dpl_rt_service/devices.d/1000.conf

sudo cp configs/dpl_rt_service.yaml /etc/kubelet.d/
```

### **Note**

The device ID and version numbers may be different in your case, please adapt as needed.

## Service Configuration

In the current release of the DPL Runtime Service, there are three types of configuration files.

Each format is similar in nature to INI files but they allow repeated sections (e.g. `[section]`) and the comments are marked by `#` rather than `;`

### Configuration of the general behavior - `dpl_rt.conf`

The DPL Runtime Service searches for this path:

```
/etc/dpl_rt_service/dpl_rt.conf
```

The contents typically look like this:

```
# Example of a possible DPL RT Service GENERAL configuration file

[LOGGING]
log_file_path=/var/log/doca/dpl_rt_service/dpl_rtd.log
log_level=INFO
# Possible log_level values (case insensitive):
# DISABLE
# CRITICAL
# ERROR
# WARNING
# INFO
# DEBUG
# TRACE

[P4RT_RPC_SERVER]
server_address=[ : : ]      # IPv6 "ANY" allows IPv4 connections
server_tcp_port=9559

[DPL_ADMIN_RPC_SERVER]
server_address=[ : : ]      # IPv6 "ANY" allows IPv4 connections
server_tcp_port=9600

[DPL_NSPECT_RPC_SERVER]
server_address=[ : : ]      # IPv6 "ANY" allows IPv4 connections
server_tcp_port=9560
```

The default logging verbosity that is configured here will be effective when the DPL Runtime Service is started. When the DPL Runtime Service is running, the logging level can be modified with the [DPL Admin client](#), provided in the DPL Dev container. Modifying the logging level from the DPL Admin tool does NOT modify the configuration file, so keep in mind that when the DPL Runtime Service is restarted, the log verbosity will be as specified in the configuration file.

This file also controls the TCP binding of three gRPC servers. It allows you to specify any address (allowing for remote connections from any accessible network interface of the

system) or to limit access to a specific IP address that is dedicated for management. Choosing a non-default TCP port for any of the gRPC servers is also possible.

## **P4RT\_RPC\_SERVER**

This is the server that listens for clients implementing the [P4Runtime protocol](#).

An [open-source client](#) is provided in the DPL Dev container.

## **DPL\_ADMIN\_RPC\_SERVER**

This is the server that listens for the p4admin client that is provided in the DPL Dev container.

## **DPL\_NSPECT\_RPC\_SERVER**

This is the server that listens to the DPL Nspect client/debugger.

## **Performance fine tuning - system.conf**

The DPL Runtime Service searches for this path:

```
/etc/dpl_rt_service/system.conf
```

```
# Example of a possible DPL RT Service system configuration file
```

```
[HAL]
queue_size=1024
queues_num=1
burst_size=32
```

These parameters control the internal behavior of how the DPL Runtime Service interfaces with the underlying hardware. Manipulating these values may affect the

maximum rate of rule insertion/deletion and its latency. Tuning these parameters for optimal values can be a complex process and is dependent on the use case and configuration of the DPU. For the current release, this file is for internal use and it is not advised for end users to change the default values, unless under the direct guidance of NVIDIA technical support.

## Device level configuration - devices.d/ .conf

The DPL Runtime Service searches for this path:

```
/etc/dpl_rt_service/devices.d/<device-id>.conf
```

e.g. `/etc/dpl_rt_service/devices.d/1000.conf`

A template is available here:

```
/etc/dpl_rt_service/devices.d/NAME.conf.template
```

In the template, you will also find internal documentation about the meaning of each field.

### Info

If you use Scalable Functions (SFs) or SR-IOV Virtual Functions (VFs), be sure to refer to their representors in the configuration file.

The `mac` and `mtu` settings are for future implementation and currently have no effect.

The values are returned back as they appear in the file when queried with the `dpl_admin` client but they should not be relied upon.

## DPL Port ID Assignment

The DPL Port IDs are assigned by the user. The user decides which DPL Port ID is assigned to which DPU interface. This mapping is critical for achieving the desired results when adding P4 table entries. A DPL Port can be:

- Uplink net device interface
- Host PF representor net device interface (`pf<X>hpf`)



- VF representor net device interface
- SF representor net device interface

### **Note**

Make sure all representors ports added to the configuration file belong to the same uplink port used.

## **DPL Port ID restrictions**

The following restrictions must be considered when assigning DPL Port IDs:

- ID of value `UINT32_MAX` is reserved
- For P4 device, the ID must be an integer number greater than zero
- For interfaces, the DPL ID must be an integer number between zero and `UINT32_MAX`
- Currently, only one Uplink port can be added to the configuration file

## **Example DPL Device Configuration File**

```
# Example of a possible DPL RT Service Device configuration file:
#
# This configuration file specifies the DPL device and its interfaces
# and their DPL Port IDs that will be used by a DPL program.
#
# The DPL Port IDs are assigned by the user. The user decides which
# DPL Port ID is assigned to which ConnectX/DPU interface. This mapping
# is critical for achieving the desired results when adding table entries.
# For DPL device, the ID must be an integer number greater than zero.
```

```

#
# The configuration file consists of following sections:
# - [DEVICE] section: Must appear only once.
# - [P4_RT_CONTROLLER] section: Must appear only once.
# - [INTERFACE] section: Must be repeated for each DPL Port (network interface).

[DEVICE]
# The DPL Device ID, used for connecting a controller to manage this device's tables.
dpl_device_id=1000
# Cache counter - decrease HW accesses - when expired an HW access will occur upon request.
dpl_counter_cache_timeout=0

[P4_RT_CONTROLLER]
# Packets delivered to the DPL RT Service from a controller will have this source DPL Port ID.
# So, this ID can be used for matching traffic originated from the controller.
p4_controller_port_id=9876

[INTERFACE]
# Interface name on the system to attach to this DPL device.
interface=p0
# DPL Port ID, used to reference this port by the DPL program and/or when updating table entries.
dpl_logical_port_id=0
# Ethernet frame size.
mtu=1514
# Only uncomment and provide this if you wish to override the interface's MAC address.
# mac=00:00:00:00:00:00

[INTERFACE]
interface=pf0hpf
dpl_logical_port_id=65535
mtu=1514
# mac=00:00:00:00:00:00

[INTERFACE]
interface=pf0vf0
dpl_logical_port_id=1
mtu=1514
# mac=00:00:00:00:00:00

```

```
[INTERFACE]
interface=pf0vf1
dpl_logical_port_id=2
mtu=1514
# mac=00:00:00:00:00:00
```

---

# DPL Development Container

The DPL Development Container is intended for developers wishing to utilize NVIDIA's DOCA Pipeline Language (DPL) to build data path applications based in a [P4](#) derived programming language. The DPL is supported on of the NVIDIA® BlueField® networking platforms (DPUs, or SuperNICs in DPU-mode).

- [DPL Target Architecture](#) is important to read for new DPL developers to understand the target architecture, its capabilities and the main building blocks that DPL applications rely upon.
- [P4 Language Support in DPL](#) outlines specific P4-16 language features supported by the DOCA Pipeline Language.
- [DPL Installation Guide](#) describes how to install the DPL Development container.
- [Compiling DPL Applications](#) describes in detail how to compile a DPL program and the various compilation options.
- [Loading DPL Applications](#) demonstrates how to load the resulting binary output of the DPL compiler on NVIDIA BlueField devices.
- [Sample DPL Applications](#) provides additional example DPL programs targeting the DPL Model.

## DOCA Target Architecture

This section introduces the NVIDIA DPL Model for the NVIDIA® BlueField®-3 networking platforms (DPUs or SuperNICs).

### Introduction

The concept of a Target Architecture (TA) is fundamental in the domain of P4 programming, serving as a [specification](#) that delineates the programmable components (such as parsers, ingress control flows, etc.) and their interfaces within a specific hardware platform. This specification acts as a formal agreement between the P4 program and the target hardware, ensuring compatibility and efficient utilization of the hardware's data

plane capabilities. For each hardware platform, vendors, including NVIDIA, provide a P4 compiler and an architecture definition that is tailored to their specific devices.

This section focuses on the DPL DOCA target architecture, offering a detailed examination of how to program these blocks to enable customized packet processing and forwarding through the various interfaces and ports available on NVIDIA's BlueField devices. Initially, the emphasis is on the BlueField-3 DPU model, with considerations for future expansion to the BlueField 4 model. While some similarities may exist between different Target Architectures, it is crucial to understand that Target Architectures are generally vendor-specific and are not designed with cross-vendor portability in mind. The following information aims to provide a comprehensive technical overview for utilizing the DPL compiler for BlueField devices, facilitating advanced data plane programming on NVIDIA's hardware platforms.

## DOCA Model

The DOCA model encompasses a set of control stages that are supported by the hardware. Each top-level object within this model is designed to receive a specific set of parameters:

- Parser – The parser is responsible for processing the input packet and a header stack. It operates without access to any metadata during the parsing stage, focusing solely on the packet's content and structure.

```
parser NvDocaParser<HEADERS>(
    packet_in packet,
    out HEADERS headers);
```

- Control – The control stage is designed to handle three distinct types of metadata, each serving a different purpose within the target architecture:
  - Standard Metadata: This is a read-only type of metadata that provides essential information about the packet obtained from prior stages, such as the parser and other fixed hardware units such as the crypto engine.
  - User Metadata: This metadata category is read/write and contains application-specific variables. It allows for the customization and storage of data relevant to the application's logic.
  - Packet Out Metadata: Also a read/write type, this metadata is associated with the packet as it is transmitted from the controller. It carries additional

information that may be required in the DOCA pipeline and for packets sent to the controller.

```
control NvDocaMainControl<HEADERS, USER_META, PKT_OUT_META>(
    inout HEADERS headers,
    in nv_standard_metadata_t std_meta,
    inout USER_META user_meta,
    inout PKT_OUT_META pkt_out_meta);
```

An important aspect of BlueField hardware is its capability to perform packet modification and reparsing as needed. Due to this feature, a deparser control stage is not necessary within this target architecture, as the native BlueField hardware can perform packet modifications and handle reparsing tasks in-line. This simplifies the main package by eliminating the need for the user to manually build a separate deparser control or recirculating the packet, streamlining packet processing within the pipeline. The main package is used as the entry point of every P4 program:

```
package NvDocaPipeline<HEADERS, USER_META, PKT_OUT_META>(
    NvDocaParser<HEADERS> parse,
    NvDocaMainControl<HEADERS, USER_META, PKT_OUT_META> main);
```

## DOCA Core Library

Basic common types are defined in `doca_core.p4`. These types are used in the DOCA TA in various interfaces such as extern functions.

### Typedefs:

The following typedefs are for commonly used header fields and standard metadata fields:

```
typedef bit<32> nv_logical_port_t;
typedef bit<48> nv_mac_addr_t;
typedef bit<32> nv_ipv4_addr_t;
```

```
typedef bit<128> nv_ipv6_addr_t;
typedef bit<32> nv_tunnel_id_t;
typedef bit<12> nv_vlan_id_t;
typedef bit<24> nv_vxlan_id_t;
typedef bit<20> nv_mpls_label_t;
typedef bit<8> nv_debug_cookie_t;
```

## Constants:

The following constants are the DOCA TA specific values for standard metadata key fields:

```
const bit<2> L2_TYPE_UNICAST    = 2w0;
const bit<2> L2_TYPE_MULTICAST  = 2w1;
const bit<2> L2_TYPE_BROADCAST  = 2w2;

const bit<2> L3_TYPE_NONE       = 2w0;
const bit<2> L3_TYPE_IPV4       = 2w1;
const bit<2> L3_TYPE_IPV6       = 2w2;

const bit<2> L4_TYPE_NONE       = 2w0;
const bit<2> L4_TYPE_TCP        = 2w1;
const bit<2> L4_TYPE_UDP        = 2w2;
const bit<2> L4_TYPE_IPSEC      = 2w3;

const bit<4> L4_TYPE_EXT_NONE   = 4w0;
const bit<4> L4_TYPE_EXT_TCP    = 4w1;
const bit<4> L4_TYPE_EXT_UDP    = 4w2;
const bit<4> L4_TYPE_EXT_ICMP   = 4w3;

const bit<2> VLAN_TYPE_NONE     = 2w0;
const bit<2> VLAN_TYPE_SVLAN   = 2w1;
const bit<2> VLAN_TYPE_CVLAN   = 2w2;

const bit<2> ENCAP_TYPE_NONE    = 2w0;
```

```

const bit<2> ENCAP_TYPE_L2_TUNNEL = 2w1;
const bit<2> ENCAP_TYPE_L3_TUNNEL = 2w2;
const bit<2> ENCAP_TYPE_ROCE      = 2w3;

const bit<2> IPSEC_TYPE_NONE      = 2w0;
const bit<2> IPSEC_TYPE_OVER_IP   = 2w1;
const bit<2> IPSEC_TYPE_OVER_UDP  = 2w2;

const bit<8> IPSEC_SYNDROME_OK    = 8w0;
const bit<8> PSP_SYNDROME_OK      = 8w0;

```

These values are used to understand the various type codes in the standard metadata.

### **#defines:**

These macros are for industry standard values for various RFC protocol fields (e.g., ethertype):

```

/* IP protocol numbers in the Protocol field of the IPv4 header
 * and the Next Header field of IPv6 header
 */
#define NV_IPV6_HBH_OPTION 0x00 /* IPv6 Hop by Hop option */
#define NV_ICMP_PROTOCOL 0x01 /* protocol number for icmp */
#define NV_IPV4_PROTOCOL 0x04 /* protocol number for IPv4 over IPv4 encap */
#define NV_TCP_PROTOCOL 0x06 /* protocol number for tcp */
#define NV_UDP_PROTOCOL 0x11 /* protocol number for udp */
#define NV_IPV6_PROTOCOL 0x29 /* protocol number for IPv6 over IPv6 encap */
#define NV_IPV6_EXT_FRAG_PROTOCOL 0x2C /* IPv6 fragmentation extension header */
#define NV_GRE_PROTOCOL 0x2F /* protocol number for Generic Routing Encapsulation */
#define NV_ESP_PROTOCOL 0x32 /* protocol number for IPsec ESP */
#define NV_AH_PROTOCOL 0x33 /* protocol number for IPsec AH */
#define NV_ICMP6_PROTOCOL 0x3A /* protocol number for icmpv6 */
#define NV_IPV6_NO_NEXT_HDR_PROTOCOL 0x3B
#define NV_SCTP_PROTOCOL 0x84 /* protocol number for SCTP */
#define NV_ROCE_PROTOCOL 0xFE /* protocol number of RoCEv1.5 */

/* Ethertype */
#define NV_TYPE_VLAN_CTAG 0x8100
#define NV_TYPE_VLAN_STAG 0x88A8

```



```

#define NV_TYPE_IPV4      0x0800
#define NV_TYPE_IPV6      0x86DD
#define NV_TYPE_ARP       0x0806
#define NV_TYPE_CONTROL   0x0808
#define NV_TYPE_MPLS      0x8847
#define NV_TYPE_MPLS_MC   0x8848
#define NV_TYPE_PTP       0x88F7
#define NV_TYPE_FCOE      0x8906
#define NV_TYPE_ROCE      0x8915 /* v1 */
#define NV_TYPE_MAC       0x6558 /* Transparent Ethernet Bridge for Generic Routing Encapsulation */

/* UDP ports */
#define NV_ROCE_PORT      4791 /* v2 bth */
#define NV_VXLAN_PORT     4789
#define NV_VXLAN_GPE_PORT 4790
#define NV_GENEVE_PORT    6081
#define NV_MPLS_TUNNEL_PORT 6635
#define NV_IPSEC_NAT_PORT 4500
#define NV_PTP_EVENT_PORT 319
#define NV_PTP_GEN_PORT   320
#define NV_IPV4_OPTION_MRI 31
#define NV_GTP_U_PORT     2152
#define NV_PSP_PORT       1000
#define NV_GUE_PORT       666

/* VXLAN GPE protocols */
#define NV_VXLAN_GPE_IPV4 0x01
#define NV_VXLAN_GPE_IPV6 0x02
#define NV_VXLAN_GPE_MAC  0x03
#define NV_VXLAN_GPE_NSH  0x04 /* RFC8300 */

/* PSP protocols */
#define NV_PSP_IPV4       0x04
#define NV_PSP_IPV6       0x29
#define NV_PSP_TCP        0x06
#define NV_PSP_UDP        0x11

```

## Metadata:

The standard metadata ( `struct nv_standard_metadata_t` ) for the DOCA TA are found in `doca_metadata.p4`. The members of the standard metadata struct are all read only, and are separated into 3 logical categories:

- Pipeline metadata, information set by hardware units before the programmable part of the BlueField pipeline
- Outer packet header metadata, information set by the hardware parser for the outer packet
- Inner packet header metadata, information set by the hardware parser for the inner packet (if present)

### User Metadata:

The DPL programmer can create and use their own metadata struct in the DOCA TA. The type can be `bit<>` or a struct, limited in size by the number of internal registers on the device. For convenience, the TA defines `nv_empty_metadata_t` which the user can pass into the main package declaration if no metadata is needed.

## DOCA Packet Parser

In the DOCA Target Architecture, headers represent the various packet protocol formats that are recognized by the parser. Once these headers are identified, the fields within the packet become candidates for further processing. They can be matched against entries in a match-action table, where they may be altered based on user-defined actions. The DOCA TA comes equipped with a default hardware-integrated parser that is capable of understanding a range of protocols standardized by the Internet Engineering Task Force (IETF). To accommodate custom requirements, the architecture allows for the integration of user-defined ("flex") headers. These flex headers are seamlessly incorporated into the pre-existing parser graph, extending its capabilities to recognize and process additional protocol formats as specified by the user. This flexibility enables developers to tailor the data plane processing to specific applications and protocols beyond the standard set, and future-proofs the developer's investment in the hardware.

This section on the DPL Packet Parser focuses on BlueField-3 devices and higher. NVIDIA's BlueField architecture combines native and flexible parsing capabilities, allowing the user to enhance the hardware's packet parsing engine with custom protocol headers. This integration is seamless, utilizing built-in parsing for standardized headers.

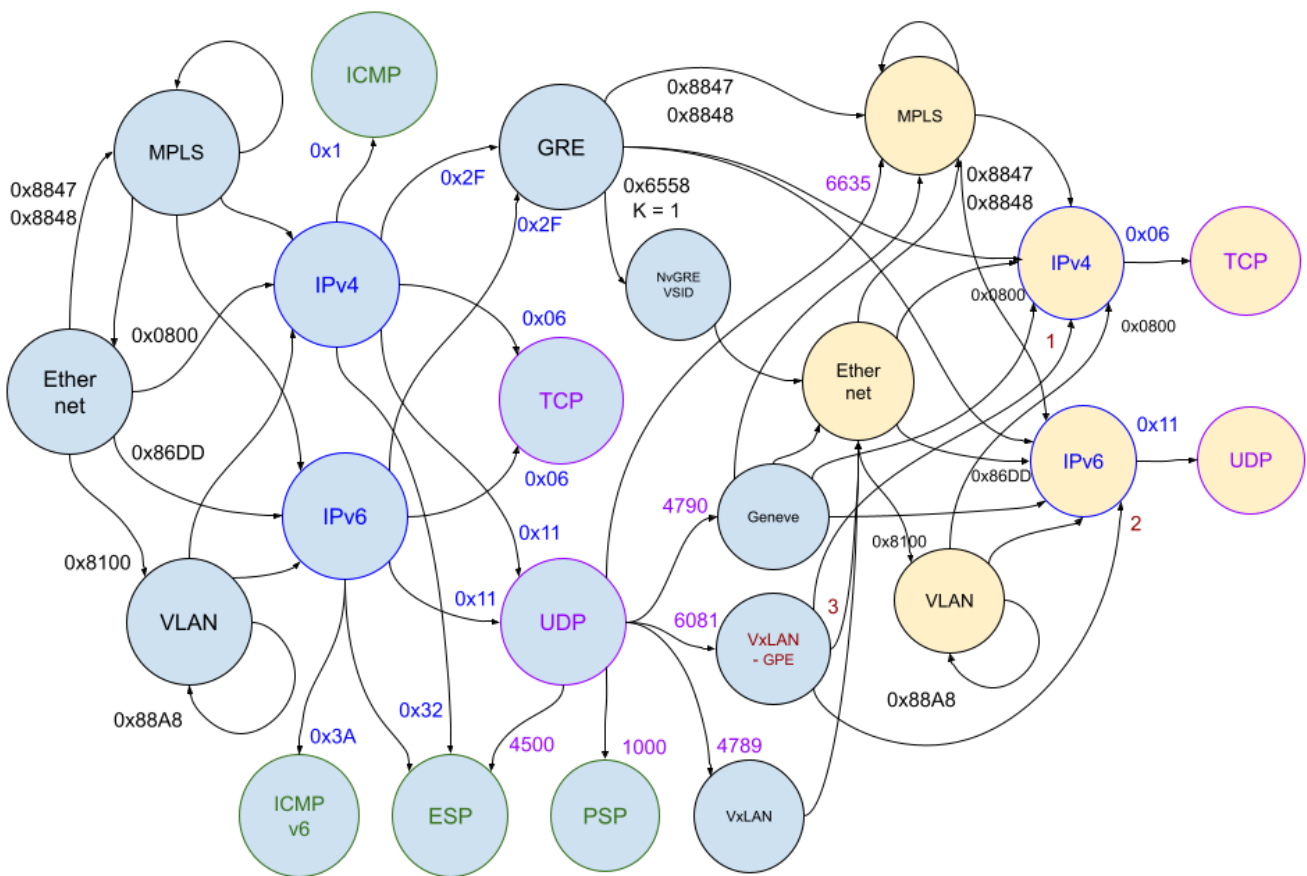
## Flex Parser Configuration

Flex nodes are user-configurable within the DOCA parser, covering elements including the next protocol field, fixed length headers, and calculations for variable header sizes. Arcs link flex nodes and feature properties such as the `is_tunnel` flag, transition value, source node, and destination node. Configurations that are not supported will trigger a compile-time error. Packet fields are extracted from the flex node to the hardware's

sampler registers, facilitating the use of fields as match keys and enabling header modification actions.

## Use Cases

Examples of custom flex headers include TCP options, SRv6, GTP-U, SCTP, GUE, eCPRI, GENEVE options, and proprietary tunnels. The hybrid architecture ensures smooth integration of user-defined protocols with the fixed parser graph. Below is an illustration of the default fixed parser for BlueField-3:



The DPL compiler enables the user to effectively utilize the P4 language to program packet parsers, taking full advantage of its native and flexible parsing capabilities to meet a variety of network processing needs. The remainder of this section describes the features of the BlueField-3 parser.

# Native Parse Graph and Headers

BlueField's parser supports a hybrid architecture that incorporates both predefined (fixed) headers and standard transitions between headers as defined by the Internet Engineering Task Force (IETF). These transitions are determined based on the "next header type" field present within each protocol. For a comprehensive understanding of these transitions and the supported protocols, users are encouraged to refer to the following resources:

- **Service Names and Port Numbers:** [IANA Service Names and Port Numbers](#)
- **IEEE 802 Numbers:** [IANA IEEE 802 Numbers](#)
- **Protocol Numbers:** [IANA Protocol Numbers](#)

These resources provide detailed information on the various service names, port numbers, IEEE standards, and protocol numbers that are recognized and supported by BlueField's parser, facilitating a broad range of packet processing applications. The following table describes the native header and transition support in the DOCA TA:

Header Type	Next Protocol Field	Field Size (bits)	Fixed IETF Transitions	Flex Transition
MAC w VLAN	<a href="#">Ether Type</a>	16	IPv4 (0x0800), IPv6 (0x86DD), MPLS UC (0x8847), MPLS MC (0x8848), VLAN (0x8100), SVLAN (0x88A8)	From, outer, inner
IPv4	<a href="#">Protocol</a>	8	UDP (0x11), TCP (0x6), GRE(0x2F), ICMP(0x1), IPsec AH (0x33), IPsec ESP (0x32), IPv4 encaps (0x4)	To, outer, inner
IPv6	<a href="#">Protocol</a>	8	HOPOPT (0x0), UDP (0x11), TCP (0x6), GRE(0x2F), ICMP(0x3A), IPsec AH (0x33), IPsec ESP (0x32), IPv6-Route (0x2B), , IPv6 encaps (0x29), IPv6-Frag(0x2C), IPv6-NoNxt (0x3B), IPv6-DestOpts (0x3c)	To, outer, inner
IP	<a href="#">Protocol</a>	8	UDP, TCP, GRE, ICMP, IPsec AH, IPsec ESP	From, outer, inner (both IPv4 and IPv6)
UDP	<a href="#">Destination port</a>	16	VXLAN (4789 + 3 additional custom ports) VXLAN-GPE (4790) GENEVE (6081) MPLS over UDP (6635)	To, from, outer, inner

Header Type	Next Protocol Field	Field Size (bits)	Fixed IETF Transitions	Flex Transition
			IPSEC ESP over UDP (4500) PSP(1000 + 2 additional custom ports)	
TCP	<a href="#">Destination port</a>	16	None	To, from, outer, inner
ICMP	None	N/A	None	Not supported - cannot transition to or from
GRE	<a href="#">Ether Type</a>	16	IPv4 (0x0800), IPv6 (0x86DD) MPLS over GRE (0x8847, 0x8848)	To, from, outer
NVGRE (GRE)	<a href="#">Ether Type</a>	16	Inner MAC (0x6558)	To, from, outer
NVGRE Options	Key present	1	Key present true	Cannot transition to or from
VXLAN	None	N/A	Inner MAC (fixed)	From, outer
VXLAN-GPE	Next Protocol	8	Reserved (0x00), IPv4 (0x01), IPv6 (0x02), Ethernet (0x03), NSH (0x04)	From, outer
GENEVE	<a href="#">Ether Type</a>	16	Inner MAC (0x6558), IPv4 (0x0800), IPv6 (0x86DD), MPLS (0x8847, 0x8848)	From, outer
MPLS	Lookahead	4	IPv4, IPv6	To, from, outer
IPSEC ESP	Next Header	8	Not supported	From, outer
<a href="#">PSP</a>	Next Header	8	Tunnel: IPv4 (4), IPv6 (41) Transport: TCP (6), UDP (17)	To, from, outer

Note that L3 control plane headers ARP (0x0806), MAC control (0x8808), LLDP (0x88CC), PTP (0x88F7) are recognized and steered to a special QP by hardware that are parsed but

not matchable in the steering pipeline.

## Using Default Native Parser

The DOCA TA defines a default native parser. The user need only include `doca_parser.p4` and reference `nv_fixed_parser` in the DPL package declaration to use the built-in hardware parser. The program should then use the native headers structure, `nv_headers_t`.

## Reparsing Capability

A distinctive feature of DPL, setting it apart from other P4 hardware architectures, is its ability to perform reparsing as needed, even midway through the processing pipeline. This flexibility allows for dynamic adjustments to packet processing based on intermediate outcomes during the pipeline's execution. Unlike RMT-based P4 Target Architectures, there are no strict pipeline stages, as the pipeline runs to completion. This means:

- Tables may be applied multiple times, saving scale by not requiring the duplication of tables and their entries at different points in the pipeline.
- Packets do not need to be resubmitted to the pipeline after packet modification, e.g. for encapsulation or decapsulation actions.
- Looping in the pipeline is permitted.

This capability greatly improves the user experience by allowing the DPL developer to focus on a logical view of the pipeline behavior, rather than dealing recirculation logic and saving state in metadata.

## Defining a Custom Parser

The P4-16 language does not provide a method to extend a pre-existing, target architecture defined parser. Hence the DPL compiler maintains an internal representation of the default parser (based on a "read only" `doca_parser.p4` source file), and allows headers to be added or removed. In this case, header removal may be implicit, based on

excluding a fixed header from the DOCA parser definition. There are 3 steps to defining a custom parser that will be used in conjunction with the native parser:

1. Define the custom header.
2. Add the header to a headers struct.
3. Define the state in a custom parser in the program.

## Custom Headers and Headers Struct

A macro, `NV_FIXED_HEADERS`, is provided that references all the headers that the fixed parser can extract. This definition may be used directly as the body of type `Headers_t`, or added to, e.g.:

```
#include <doca_parser.p4>
header custom_header_t {
    bit<16> x;
}
struct custom_headers_t
{
    NV_FIXED_HEADERS
    custom_header_t custom;
}
```

The DPL programmer is expected to use `NV_FIXED_HEADERS` as the basis of their `Headers_t` definition. Failure to do so may result in unexpected behavior, for example:

- Headers/fields that are never extracted at runtime even if extracted by parser states.
- Use of flex extractions for every single referenced header field.
- Compilation errors.

## Custom Parser States

A macro, `NV_FIXED_PARSER`, is provided which describes the fixed parser. This definition may be used directly as the body of a parser, or extended by writing additional states, and linking them using the `@nv_transition_from` annotation:

```
parser my_parser(packet_in packet, out custom_headers_t headers)
{
    NV_FIXED_PARSER(packet, headers)

    @nv_transition_from("nv_parse_ethernet", 0x1234)
    state parse_custom {
        packet.extract(headers.custom);
        transition select(headers.custom.x) {
            1: nv_parse_ipv4;
            2: nv_parse_ipv6;
            default: accept;
        }
    }
}
```

The DPL programmer is expected to use `NV_FIXED_HEADERS` as the basis of their parser definition. The macro must take the same variable name use for `packet_in` and `headers` used in the parser control definition. Failure to do so may result in unexpected behavior, for example:

- Use of hardware flex parser resources for every parser state, even where hardware fixed parser resources could have been used.
- Definitions of states with names matching states in `NV_FIXED_PARSER` being ignored, and their fixed definition used instead.
- Compilation errors.

Finally, the custom parser should be instantiated in the main DPL package:

```
NvDocaPipeline(
    my_parser(),
```



```
    my_main_control()  
  ) main;
```

## Specifying Transitions to Flex Headers

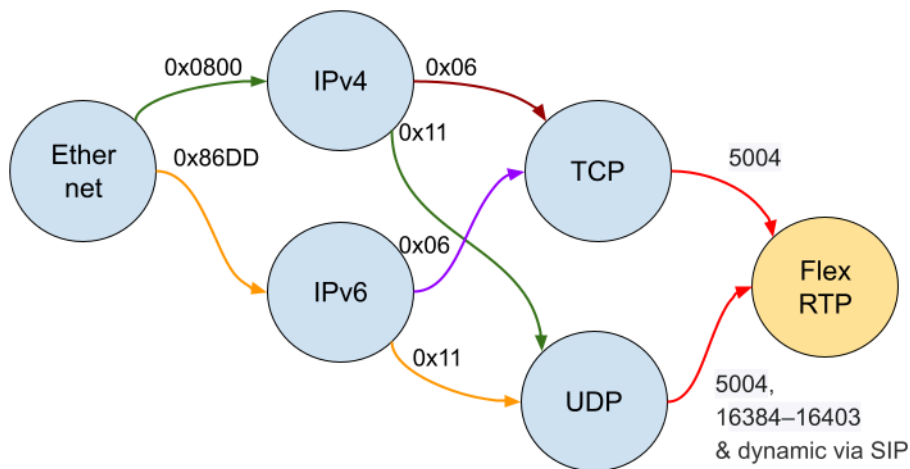
To extend the native parser with flex headers, the DOCA TA utilizes a combination of macros and DPL annotations to instruct the compiler how to stitch in the customized headers into an existing parse graph.

To create an unconditional transition from a source state, the parser state must be preceded by an `@nv_transition_from` annotation with 1 argument: the source state name. If the parser contains multiple state definitions with an `@nv_transition_from` annotation with the same source state, a single state will be selected (in an undefined manner) to be the target state of the transition. Put another way, each `@nv_transition_from` annotation from a given source state annotation may be considered to override any earlier (in an undefined order) transition destination for that source state.

To create a conditional transition based on the source state's transition field, the parser state must be preceded by an `@nv_transition_from` annotation with 2 arguments: the source state name and transition value. Transitions of this type are prepended to the source state's list of transitions and thus are guaranteed to override any transitions defined by `NV_FIXED_PARSER` in an undefined order. Entries with identical transition values will override each other but may generate warnings due to unreachability.

The following diagram shows an example of extending the native parser to support the RTP protocol:

## RFC 1889, 3550, 3551: RTP



### **i** Note

In the native parser, as an optimization, the user can specify a transition from both IPv4 and IPv6 headers to a flex header using a special 2 argument annotation of the format:

```
@nv_transition_from("nv_parse_ipv6, nv_parse_ipv4", 0x1234)
```

Where `0x1234` is an example of the custom IP protocol value for some user defined L4 protocol.

## Specifying Transitions to Fixed Headers

In a limited number of cases, transitions from fixed headers can be configured. This uses a different annotation, `@nv_transition`, and is always used to annotate the top level parser object. For example, it is a common requirement to add nonstandard transitions from UDP to VxLAN:

```

@nv_transition("nv_parse_udp", "nv_parse_vxlan", 1234)
@nv_transition("nv_parse_udp", "nv_parse_vxlan", 1235)
@nv_transition("nv_parse_udp", "nv_parse_vxlan", 1236)
parser custom_parser(
    packet_in packet,
    out nv_headers_t headers
) {
    NV_FIXED_PARSER(packet, headers)
}

```

where the first parameter specifies the parser node to transition "from", and the second parameter specifies the parser node to transition "to". Note that the annotation `@nv_transition` can be used instead of `@nv_transition_from` for connecting headers.

## Defining TLV Parsing for Headers with Optional Data

The DOCA TA allows the user to define protocol headers that have TLVs as sub headers. Note that a custom base "parent" header must be supplied in order to create custom TLVs. It is not supported to create custom TLVs for native headers. For example, the user can define a GENEVE header as follows:

```

header geneve_t {
    bit<2> ver;
    bit<6> opt_len;
    bit<1> o;
    bit<1> c;
    bit<6> reserved;
    bit<16> protocol_type;
    bit<24> vni;
    bit<8> reserved56;
};

```

And then define a struct for the base GENEVE option fields (TYPE and LENGTH) along with a customized VALUE data.

```
// Struct so it can be a field within other headers
struct geneve_option_t {
    bit<24> option_class;
    bit<8> option_type;
    bit<3> reserved24;
    bit<5> length;
};

header geneve_option_int_md_t {
    geneve_option_t base;
    bit<4> ver;
    bit<1> d;
    bit<1> e;
    bit<1> m;
    bit<12> reserved7;
    bit<5> hop_m1;
    bit<8> remaining_hop_count;
    // etc.
};
```

Next, the user must add the custom headers to the struct of headers and connect them in the parser. Both the base option struct and the newly defined headers must be present in the headers struct so that there is a reference to these types in the parser.

```
struct app_headers {
    NV_FIXED_HEADERS
    geneve_t custom_geneve;
    geneve_option_t geneve_opt;
    geneve_option_int_md_t geneve_opt_int_md;
};
```

The DOCA TA uses an `NvOptionsParser` extern object to define TLV options, since the P4 language itself does not provide a native way to define TLV parsing. An instance of the `NvOptionsParser` must be created in the parser, at the outermost scope, with the following parameters:

- `options_length_field` - the name of the field in the parent header (as a string) that holds the total length of the options. If the total options length is fixed (i.e., not specified in the header's data), then this field should be omitted.
- `options_length_shift` - the value by which the parser should apply a shift (i.e. multiply by power of 2) to the value of the option length field. If the length value is fixed or does not need to be shifted, then this field should be omitted.
- `options_length_add` - the constant to which the parser should add to the value option length field after shifting. If the length value does not need an addend, then this field should be omitted. If the total options length is fixed, then this field must specify the fixed length of the TLV options.
- `option_layout_header_type` - the user defined P4 struct that holds the TLV field layout of the option header
- `option_length_field` - the name of the field in the base option struct (as a string) that holds the length of the options. If the options length is fixed (i.e., not specified in the base options struct), then this field should be omitted.
- `option_length_shift` - the value by which the parser should apply a shift (i.e. multiply by power of 2) to the option length field. If the length value is fixed or does not need to be shifted, then this field should be omitted.
- `option_length_add` - the constant to which the parser should add to the option length field. If the length value does not need an addend, then this field should be omitted. If the options length is fixed, then this field must specify the fixed length of the TLV options.
- `options` - a list of tuples, where the first value is the TYPE, and the second value is the name of the child header defined by the user

```
parser geneve_parser(  
    packet_in packet,  
    out app_headers headers
```

```

) {
    NvOptionParser<bit<24>, _>(
        "opt_len",          // options_length_field
        2,                  // options_length_shift
        0,                  // options_length_add
        "geneve_option_t", // option_layout_header_type
        "length",          // option_length_field
        0,                  // option_length_shift
        4,                  // option_length_add
        "option_class_type", // option_type_field
        // options data
        (list<tuple<bit<24>, _>>){
            {24w0x010301, "headers.geneve_opt_int_md"}
            // list of additional options ...
        }
    ) geneveOptions;

    NV_FIXED_PARSER(packet, headers)

    @nv_transition_from("nv_parse_udp", 6082)
    state parse_custom_geneve {
        packet.extract(headers.custom_geneve);
        geneveOptions.parseOptions(packet, headers);
        transition accept;
    }
}

```

## Restrictions and Unsupported Parser Features

DOCA parser features that either differ from the P4<sub>16</sub> specification, are restricted in the DOCA TA or are unsupported are listed below.

### Unsupported Features

These features are not supported by the DOCA TA in comparison to the P4<sub>16</sub> specification. Utilizing these features will result in compilation errors.

- Variable Declarations and Extern Instantiations: Parsers and parser states cannot contain any variable declarations or extern instantiations other than `NvOptionsParser`.
- Lookahead and Advance: Programmable lookahead and advance functionalities are not supported.
- Parser Value Sets: Parser value sets for setting transition values from runtime data are not supported.
- State Declarations: States may not contain any declaration, or any statement besides extract, transition statements and calls to `NvOptionsParser` methods.

## Restricted Behavior

The following points below highlight how the DOCA TA implementation of DPL parser features differs from the standard P4<sub>16</sub> specification.

- Reserved State Names: Parser state names beginning with `nv_*` are reserved by the DOCA TA. Custom states using this reserved prefix will be ignored.
- State and Header Coupling: Each state must extract exactly one header, and they are considered to be coupled in a 1:1 fashion.
- Fixed-Size Header Extraction: The extracted header must be fixed-size, known at compile-time, and defined in the P4 header type definition. I.e., `p.extract(headers.xxx)` is allowed, but `p.extract(headers.xxx, someLength)` is not.
- Transition Statements: Besides terminal pre-defined states (`accept`, `reject`), states must transition to another state using a transition statement, which may be unconditional or conditional.
- Transition Select: Transition select statements are limited to using a single field from the header the state extracts or a constant. Lists of expressions, constants, operators, references to declarations other than a field in the header the state extracts, and any other expression not explicitly allowed are disallowed.

- Default Case in Transition Select: Every transition select statement must have a default case, transitioning to the accept state.
- Extract Before Transition: A state's `extract` statement must precede its transition statement.
- Empty Accept and Reject States: The `accept` and `reject` states must be completely empty.
- Loops in Parser Graph: Loops are not allowed in the DPL program's parser graph. Each path through the parser must be acyclic.

## Hybrid Parser Behavior

These restrictions and behaviors apply to the unique hybrid fixed/flex parser model supported by the DOCA TA.

- Start state of `nv_parse_ethernet` – The start state cannot be changed or overwritten.
- Redefining states – States may not be redefined. This includes no possibility of applying `@nv_transition_from` to the body of any fixed state.
- `NV_FIXED_PARSER` and `NV_FIXED_HEADERS` – States part of `NV_FIXED_PARSER` extract the pre-defined fixed header, irrespective of the state body. States not part of `NV_FIXED_PARSER` must extract headers not part of `NV_FIXED_HEADERS`.
- Fixed-to-fixed transitions – It is not possible to disable any fixed->fixed transition. These may be overridden to transition to a custom state by applying an `@nv_transition_from` annotation to the destination state.
- MPLSoUDP stops parsing at the inner MAC header if the pseudowire Control Word is set to 0. Additional inner packet fields will not be parsed.

## DOCA Extern Objects



# Counters

BlueField supports per entry counters, shared and direct. A (shared) counter is an extern object that allows the user to access N independent counters via an index. Whereas a direct counter is an extern object that is directly associated with a P4 table where the index is implied by the entry slot. For both types, the counter value is updated when an action calls the count() method.

## NvCounter

Instantiates an indirect counter of the specified width and type.

**Signature:** `extern NvCounter(bit<32> size, NvCounterType type)`

### Parameters:

- `size[in]`: Number of counter indices
- `type[in]`: Enum indicating counter type. Currently only `NvCounterType.PACKETS_AND_BYTES` is supported.

### Methods:

#### count

Increments the counter at the specified index

```
void count(in bit<32> index)
```

### Parameters:

- `index[in]`: Index of the counter to increment

## NvDirectCounter

Instantiates a direct counter that can be assigned to a table. Direct counters do not have a separate size, they are as wide as the table. Each action invoked by an entry hit should call the count() method to increment the entry count.

**Signature:** `extern NvDirectCounter(NvCounterType type)`

## Parameters:

- `type[in]`: Enum indicating counter type. Currently only `NvCounterType.PACKETS_AND_BYTES` is supported.

## Methods:

### count

Increments the counter associated the the entry that triggered the action execution

```
void count()
```

#### Note

Methods of a `NvDirectCounter` object may only be called by the P4 Table that owns it via the "direct\_counter" property bound to the P4 Table. It is an error to attempt to call the `count()` method from a P4 Table that does not own the direct counter.

#### Note

If there are no side effects, it is most efficient to call the `count` method as the first primitive action in a P4 action.

## Meters

BlueField supports per-flow metering, following RFC2697, RFC2698, and RFC4115. Similar to counters, both direct meters and shared meters are supported in the DPL target architecture. The packets for a given flow can be colored by the meter into three "colors": red, yellow and green. The meter burst parameters can be specified with units of bytes or packets.

```

typedef bit<8> nv_meter_color_t;

enum nv_meter_color_t NvMeterColor {
    RED = 8w0,
    YELLOW = 8w1,
    GREEN = 8w2,
}

enum NvMeterUnits {
    BYTES,
    PACKETS
}

```

### **Note**

Each configuration parameter of NvMeterPeakTrTCM and NvDirectMeterPeakTrTCM is restricted to certain max values.

When metering by BYTES, max values are the following:

- cir, pir: 255000000000
- cbs, pbs: 2147483648

When metering by PACKETS, max values are the following:

- cir, pir: 1992187500
- cbs, pbs: 16777216

## **NvMeterPeakTrTCM**

Instantiates a Two Rate Three Color Marker shared meter object of the specified size, units and burst parameters. (Implements [RFC 2698](#))

**Signature:**

```
extern NvMeterPeakTrTCM(bit<32> size, NvMeterUnits units, bit<64>
cir, bit<64> cbs, bit<64> pir, bit<64> pbs)
```

**Parameters:**

- `size[in]`: Number of meter indices
- `units[in]`: Enum indicating the units used by the meter. Metering by `NvMeterUnits.BYTES` or by `NvMeterUnits.PACKETS` is supported.
- `cir[in]`: Committed Information Rate
- `cbs[in]`: Committed Burst Size
- `pir[in]`: Peak Information Rate
- `pbs[in]`: Peak Burst Size

**Methods:****meter**

Executes the metering of the flow entry at the specified index

```
nv_meter_color_t meter(in bit<32> index, nv_meter_color_t
initial_color = NvMeterColor.GREEN)
```

**Parameters:**

- `index[in]`: Index of the meter to execute
- `initial_color[in]`: The initial color to assign to the flow. Coloring of `NvMeterColor.RED`, `NvMeterColor.YELLLOW`, or `NvMeterColor.GREEN`. The default is `NvMeterColor.GREEN`.

**Returns:**

- The `nv_meter_color_t` of the flow entry at the specified index.

## NvDirectMeterPeakTrTCM

Instantiates a Single Rate Three Color Marker indirect meter object of the specified size, units and burst parameters. (Implements [RFC 2698](#))

### Signature:

```
extern NvDirectMeterPeakTrTCM(bit<32> size, NvMeterUnits units,  
bit<64> cir, bit<64> cbs, bit<64> pir, bit<64> pbs)
```

### Parameters:

- `size[in]`: Number of meter indices
- `units[in]`: Enum indicating the units used by the meter. Metering by `NvMeterUnits.BYTES` or by `NvMeterUnits.PACKETS` is supported.
- `cir[in]`: Committed Information Rate
- `cbs[in]`: Committed Burst Size
- `pir[in]`: Peak Information Rate
- `pbs[in]`: Peak Burst Size

### Methods:

#### meter

Executes the metering of the flow entry at the specified index

```
nv_meter_color_t meter(nv_meter_color_t initial_color =  
NvMeterColor.GREEN)
```

### Parameters:

- `initial_color[in]`: The initial color to assign to the flow. Coloring of `NvMeterColor.RED`, `NvMeterColor.YELLLOW`, or `NvMeterColor.GREEN`. The default is `NvMeterColor.GREEN`.

### Returns:

- The `nv_meter_color_t` of the flow entry at the specified index.

### **Note**

Methods of a **NvDirectMeterPeakTrTCM** object may only be called by the P4 Table that owns it via the "direct\_meter" property bound to the P4 Table. It is an error to attempt to call the meter() method from a P4 Table that does not own the direct meter. A meter may not be owned by multiple P4 Tables.

## Options Parser

DPL supports TLV style parsing for protocols with options.

### NvOptionsParser

The NvOptionsParser is an extern object used to define the behavior of the TLV (Type-Length-Value) options parser.

```
extern NvOptionParser<VALUE, OPTION> {
    NvOptionParser(
        string options_length_field = "",
        bit<32> options_length_shift = 32w0,
        bit<32> options_length_add = 32w0,
        string option_layout_header_type = "",
        string option_length_field = "",
        bit<32> option_length_shift = 32w0,
        bit<32> option_length_add = 32w0,
        string option_type_field = "",
        list<tuple<VALUE, OPTION>> options
    );
};
```

```
void parseOptions<H>(packet_in p, inout H headers);  
}
```

### Parameters:

- options\_length\_field[in]: Optional; omit if **total** options length is fixed, otherwise provide the field name that contains the header length value
- options\_length\_shift: Optional; omit if the field can be used without shifting
- options\_length\_add: Optional; omit if the field can be used without adding
- option\_layout\_header\_type: Mandatory header\_type\_name
- option\_length\_field: Optional; omit if **each** option is a constant size, otherwise provide the field name that contains the header length value
- option\_length\_shift: Optional; omit if the field can be used without shifting
- option\_length\_add: Optional; omit if the field can be used without adding
- option\_type\_field: Mandatory field\_name
- options: Mandatory; variable-length (expression-)list of (type, header field)

### Methods:

#### parseOptions

Parses the TLV

**Signature:** `void parseOptions<H>(packet_in packet, inout H headers);`

### Parameters:

- packet[in]: input packet
- headers[in/out]: packet headers struct

### Examples:

See [Geneve TLV Parsing Example](#).

# Custom Tunnel Encapsulation

DPL supports the creation of L2 and L3 tunnels using custom tunnel headers.

## NvTunnelTemplate

Instantiates a template object used to create tunnel encapsulations. This object is used for controlling the behavior of the actions `nv_set_l2tunnel_underlay` and `nv_set_l3tunnel_underlay`.

```
extern NvTunnelTemplate<HEADER_TYPE> {  
    NvTunnelTemplate();  
}
```

### Type Variables:

- `HEADER_TYPE[in]`: Must be the typename of a user-declared `struct` whose only fields must be `header` types. The entire set of underlay headers should be defined as a struct, for example:

```
struct tunnel_headers_t {  
    nv_ethernet_h ethernet;  
    my_tunnel_h    custom_tunnel;  
}
```

### Annotation:

NvTunnelTemplate works slightly differently than other extern objects in that the object has no methods; instead, the object can be passed into one of two methods `nv_set_l2tunnel_underlay` and `nv_set_l3tunnel_underlay`. Declaring a NvTunnelTemplate object does not use any resources in hardware, as the template construct is provided for simplifying the process of specifying header field values for custom tunnels. When declaring a NvTunnelTemplate object, the annotation



`@nv_tunnel_fields` must be present. This annotation contains a key-value entry for each header in the specified struct type. For each header, the header fields must be specified in the order that they appear in the header definition. The value assigned to a field must be one of the following:

- a non-negative integer
- the string "variable"
- the string "ignore"

```
@nv_tunnel_fields(  
    ethernet = {  
        dst_addr = "variable",  
        src_addr = "variable",  
        ether_type = 0xABCD  
    },  
    custom_tunnel = {  
        field_1 = 1,  
        field_2 = "ignore",  
        field_3 = "variable"  
    }  
)  
NvTunnelTemplate<tunnel_headers_t>() my_tunnel;
```

For a complete example, see [GTP Tunnel Encapsulation Example](#).

### **Note**

Certain header fields will automatically be recalculated by the hardware after encapsulation. These fields cannot be assigned by the user. It is good practice to mark them as "ignore".

Header field	Value
ipv4.ecn	0

Header field	Value
ipv4.identification	0
ipv4.hdr_checksum	Hardware calculated
ipv4.total_len	Hardware calculated
ipv6.traffic_class (ECN bits only)	0
ipv6.payload_length	Hardware calculated
udp.src_port	Entropy hash
udp.length	Hardware calculated
udp.checksum	0
tcp.checksum	Hardware calculated
gre.key (8 LSB only)	Entropy hash

## DOCA Extern Functions

Extern functions serve as a mechanism for exposing DPL target specific functionality that may be beyond a standard P4 model.

Note that many of these externs mutate the packet, for which the HEADERS are a part of the extern signature.

### nv\_drop

Terminal extern function that stops packet processing and drops the packet

**Signature:** `extern void nv_drop();`

### nv\_send\_to\_port

Terminal extern function that stops packet processing and sends the packet to the specified port

**Signature:** `extern void nv_send_to_port(in nv_logical_port_t port);`

**Parameters:**

- port[in]: This parameter specifies the logical port to which the packet will be sent to

## nv\_send\_to\_controller

Terminal extern function that forwards packet metadata to a controller

**Signature:**

```
extern void nv_send_to_controller<PACKET_META>(in PACKET_META
pkt_in_meta);
```

**Parameters:**

- pkt\_in\_meta[in]: The metadata of the packet being sent to the controller for processing

## nv\_dec\_ip\_ttl

Extern function that decrements the ttl value in the IP header. Applies to both IPv4 and IPv6. The ttl value will be set to max value if decremented from zero value.

**Signature:**

```
extern void nv_dec_ip_ttl<HEADERS>(inout HEADERS headers, in bit<8>
ttl_value);
```

**Parameters:**

- headers[in/out]: packet headers struct
- ttl\_value[in]: TTL value to decrement

## nv\_set\_ip\_dscp

Extern function that sets the DSCP value in the IP header. Applies to both IPv4 and IPv6.

**Signature:**

```
extern void nv_set_ip_dscp<HEADERS>(inout HEADERS headers, in bit<6>
dscp_value);
```

**Parameters:**

- headers[in/out]: packet headers struct
- dscp\_value[in]: DSCP value to set

## nv\_set\_ip\_ecn

Extern function that sets the ECN value in the IP header. Applies to both IPv4 and IPv6.

**Signature:**

```
extern void nv_set_ip_ecn<HEADERS>(inout HEADERS headers, in bit<2>
ecn_value);
```

**Parameters:**

- headers[in/out]: packet headers struct
- ecn\_value[in]: ECN value to set

## nv\_set\_ip\_ttl

Extern function that sets the TTL value in the IP header. Applies to both IPv4 TTL and IPv6 hop limit.

**Signature:**

```
extern void nv_set_ip_ttl<HEADERS>(inout HEADERS headers, in bit<8>
ttl_value);
```

**Parameters:**

- headers[in/out]: packet headers struct

- ttl\_value[in]: TTL/Hop limit value to set

## nv\_set\_l4\_src\_port

Extern function that sets the L4 source port value in the TCP/UDP header.

### Signature:

```
extern void nv_set_l4_src_port<HEADERS>(inout HEADERS headers, in  
bit<16> src_port);
```

### Parameters:

- headers[in/out]: packet headers struct
- src\_port[in]: L4 port value to set

## nv\_set\_l4\_dst\_port

Extern function that sets the L4 destination port value the TCP/UDP header.

### Signature:

```
extern void nv_set_l4_dst_port<HEADERS>(inout HEADERS headers, in  
bit<16> dst_port);
```

### Parameters:

- headers[in/out]: packet headers struct
- dst\_port[in]: L4 port value to set

## nv\_l2\_decap

Extern function that performs Layer 2 decapsulation on the packet headers i.e. removing ethernet/vlan headers.

Signature: `void nv_l2_decap<HEADERS>(inout HEADERS headers)`

**Parameters:**

- headers[in/out]: packet headers struct that will be decapsulated

## nv\_l3\_decap

Extern function that performs Layer 3 decapsulation on the packet headers i.e. removing ethernet, vlan and IP headers, and appends new L2 headers with specified.

```
void nv_l3_decap<HEADERS>(inout HEADERS headers, bool has_vlan,  
nv_mac_addr_t dst_mac, nv_mac_addr_t src_mac, bit<16> l3_ether_type,  
nv_vlan_id_t vid)
```

**Parameters:**

- headers [inout]: The packet headers that will be decapsulated
- has\_vlan [in]: Bool flag indicating whether the packet includes a VLAN header.
- dst\_mac [in]: The destination MAC address in the Ethernet frame.
- src\_mac [in]: The source MAC address in the Ethernet frame.
- l3\_ether\_type [in]: The EtherType value in the Ethernet frame
- vid [in]: The VLAN ID, only valid if has\_vlan is true.

## nv\_push\_vlan

Extern function that inserts a VLAN header immediately after the L2 header.

```
void nv_push_vlan<HEADERS>(inout HEADERS headers, in NvVlanTagId  
tpid, in bit<3> pcp, in bit dei, in nv_vlan_id_t vid);
```

**Parameters:**

- headers [inout]: The packet headers that will be decapsulated.
- tpid [in]: Tag Protocol Identifier.

- pcp [in]: Priority Code Point.
- dei [in]: Drop Eligible Indicator (formerly CFI).
- vid [in]: VLAN Identifier.

## nv\_pop\_vlan

Extern function that removes the outermost VLAN header, immediately after the L2 header.

```
void nv_pop_vlan<HEADERS>(inout HEADERS headers);
```

### Parameters:

- headers [inout]: The packet headers that will be decapsulated.

## nv\_set\_vxlan\_v4\_underlay

### Description:

Extern function that encapsulates the packet with an ethernet frame (optionally VLAN tagged), an ipv4 header and a VXLAN header.

### Signature:

```
void nv_set_vxlan_v4_underlay<HEADERS>(inout HEADERS headers, in
bool has_vlan, in nv_mac_addr_t dst_mac, in nv_mac_addr_t src_mac,
in nv_vlan_id_t vid, in nv_ipv4_addr_t sip, in nv_ipv4_addr_t dip,
in nv_vxlan_id_t vni)
```

### Parameters:

- headers [in/out]: The packet headers that are to be encapsulated with a VXLAN underlay.
- has\_vlan [in]: Specifies whether a VLAN tag is present.
- dst\_mac [in]: Destination MAC address.

- `src_mac [in]`: Source MAC address.
- `vid [in]`: VLAN ID (only valid if `has_vlan` is set, otherwise ignored)
- `sip [in]`: Source IPv4 address.
- `dip [in]`: Destination IPv4 address.
- `vni [in]`: VXLAN Network Identifier.

## nv\_set\_vxlan\_v6\_underlay

### Description:

Extern function that encapsulates the packet with an ethernet frame (optionally VLAN tagged), an IPv6 header, and a VXLAN header.

### Signature:

```
void nv_set_vxlan_v6_underlay<HEADERS>(inout HEADERS headers, in
bool has_vlan, in nv_mac_addr_t dst_mac, in nv_mac_addr_t src_mac,
in nv_vlan_id_t vid, in nv_ipv6_addr_t sip, in nv_ipv6_addr_t dip,
in nv_vxlan_id_t vni)
```

### Parameters:

- `headers [in/out]`: The packet headers that are to be encapsulated with a VXLAN underlay.
- `has_vlan [in]`: Specifies whether a VLAN tag is present.
- `dst_mac [in]`: Destination MAC address.
- `src_mac [in]`: Source MAC address.
- `vid [in]`: VLAN ID (only valid if `has_vlan` is set, otherwise ignored).
- `sip [in]`: Source IPv6 address.
- `dip [in]`: Destination IPv6 address.



- `vni` [in]: VXLAN Network Identifier.

## nv\_set\_gre\_v4\_underlay

### Description:

Extern function that configures GRE IPV4 underlay encapsulation, with an optional VLAN tag and optional GRE key.

### Signature:

```
void nv_set_gre_v4_underlay<HEADERS>(inout HEADERS headers, in bool has_vlan, in bool has_key, in nv_mac_addr_t dst_mac, in nv_mac_addr_t src_mac, in nv_vlan_id_t vid, in nv_ipv4_addr_t sip, in nv_ipv4_addr_t dip, in NvInnerProtocolType proto_type, in bit<32> key)
```

### Parameters:

- `headers` [in/out]: The packet headers that are to be encapsulated with a GRE underlay.
- `has_vlan` [in]: Specifies whether a VLAN tag is present.
- `has_key` [in]: Specifies whether the GRE key is present.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set, otherwise ignored).
- `sip` [in]: Source IPv4 address.
- `dip` [in]: Destination IPv4 address.
- `proto_type` [in]: Protocol type for the inner payload.
- `key` [in]: GRE key (only valid if `has_key` is set, otherwise ignored).

## nv\_set\_gre\_v6\_underlay

### Description:

Extern function that configures GRE IPV6 underlay encapsulation, with an optional VLAN tag and optional GRE key.

### Signature:

```
void nv_set_gre_v6_underlay<HEADERS>(inout HEADERS headers, in bool has_vlan, in bool has_key, in nv_mac_addr_t dst_mac, in nv_mac_addr_t src_mac, in nv_vlan_id_t vid, in nv_ipv6_addr_t sip, in nv_ipv6_addr_t dip, in NvInnerProtocolType proto_type, in bit<32> key)
```

### Parameters:

- `headers` [in/out]: The packet headers that are to be encapsulated with a GRE underlay.
- `has_vlan` [in]: Specifies whether a VLAN tag is present.
- `has_key` [in]: Specifies whether the GRE key is present.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set, otherwise ignored).
- `sip` [in]: Source IPv6 address.
- `dip` [in]: Destination IPv6 address.
- `proto_type` [in]: Protocol type for the inner payload.
- `key` [in]: GRE key (only valid if `has_key` is set, otherwise ignored).

# nv\_set\_geneve\_v4\_underlay

## Description:

Extern function that configures GENEVE IPV4 underlay encapsulation with optional VLAN tagging and an optional GENEVE option with 32 bits of data.

## Signature:

```
void nv_set_geneve_v4_underlay<HEADERS>(inout HEADERS headers, in
bool has_vlan, in bool has_option, in nv_mac_addr_t dst_mac, in
nv_mac_addr_t src_mac, in nv_vlan_id_t vid, in nv_ipv4_addr_t
src_ip, in nv_ipv4_addr_t dst_ip, in bit<6> opt_fields_len, in
bit<1> oam, in bit<1> critical, in NvInnerProtocolType proto_type,
in bit<24> vni, in bit<16> opt_class, in bit<8> opt_type, in bit<5>
opt_len, in bit<32> opt_data)
```

## Parameters:

- `headers` [in/out]: The packet headers that are to be encapsulated with a GENEVE underlay.
- `has_vlan` [in]: Specifies whether a VLAN tag is present.
- `has_option` [in]: Specifies whether a GENEVE option is present.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set, otherwise ignored).
- `src_ip` [in]: Source IPv4 address.
- `dst_ip` [in]: Destination IPv4 address.
- `opt_fields_len` [in]: Length of the options field in 4-byte units (must be set to 0 if `has_option` is false)
- `oam` [in]: OAM flag.

- `critical` [in]: Critical options flag.
- `proto_type` [in]: Protocol type for the inner payload.
- `vni` [in]: Virtual Network Identifier.
- `opt_class` [in]: Option class (ignored if `has_option` is false)
- `opt_type` [in]: Option type. (ignored if `has_option` is false)
- `opt_len` [in]: Option length in 4-byte units. (ignored if `has_option` is false, otherwise must be 1)
- `opt_data` [in]: Option data. (ignored if `has_option` is false)

## nv\_set\_geneve\_v6\_underlay

### Description:

Extern function that configures GENEVE IPV6 underlay encapsulation with optional VLAN tagging and an optional GENEVE option with 32 bits of data.

### Signature:

```
void nv_set_geneve_v6_underlay<HEADERS>(inout HEADERS headers, in
bool has_vlan, in bool has_option, in nv_mac_addr_t dst_mac, in
nv_mac_addr_t src_mac, in nv_vlan_id_t vid, in nv_ipv6_addr_t
src_ip, in nv_ipv6_addr_t dst_ip, in bit<6> opt_fields_len, in
bit<1> oam, in bit<1> critical, in NvInnerProtocolType proto_type,
in bit<24> vni, in bit<16> opt_class, in bit<8> opt_type, in bit<5>
opt_len, in bit<32> opt_data)
```

### Parameters:

- `headers` [in/out]: The packet headers that are to be encapsulated with a GENEVE underlay.
- `has_vlan` [in]: Specifies whether a VLAN tag is present.

- `has_option` [in]: Specifies whether a GENEVE option is present.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set, otherwise ignored).
- `src_ip` [in]: Source IPv6 address.
- `dst_ip` [in]: Destination IPv6 address.
- `opt_fields_len` [in]: Length of the options field in 4-byte units (must be set to 0 if `has_option` is false, otherwise must be set to 2)
- `oam` [in]: OAM flag.
- `critical` [in]: Critical options flag.
- `proto_type` [in]: Protocol type for the inner payload.
- `vni` [in]: Virtual Network Identifier.
- `opt_class` [in]: Option class (ignored if `has_option` is false)
- `opt_type` [in]: Option type. (ignored if `has_option` is false)
- `opt_len` [in]: Option length in 4-byte units. (ignored if `has_option` is false, otherwise must be 1)
- `opt_data` [in]: Option data. (ignored if `has_option` is false)

## **nv\_set\_l2tunnel\_underlay**

### **Description:**

Extern function that configures custom underlay followed by an L2 encapsulation of the original packet. It is optional, but recommended, that the new tunnel encapsulation be parsable by the DPL parser. For more details see the [section on NvTunnelTemplate](#).

## Signature:

```
void nv_set_l2tunnel_underlay<HEADERS, NV_TUNNEL_TEMPLATE, VALUES>(
  inout HEADERS headers, NV_TUNNEL_TEMPLATE tunnel, in VALUES values)
```

## Parameters:

- `headers` [in/out]: The packet headers that are to be encapsulated.
- `tunnel` [in]: The `NvTunnelTemplate` object that specifies the underlay headers and the header values.
- `values` [in]: A list expression whose length matches the number of header fields marked "variable" in the `NvTunnelTemplate` object's annotation. The order of the header fields determines which value a field corresponds to, i.e. 3rd header field marked "variable" will be set to the 3rd value in the list. Each value is required to be runtime constant, meaning it may either be an integer value or an action parameter value.

## nv\_set\_l3tunnel\_underlay

### Description:

Extern function that configures custom underlay followed by an L3 encapsulation of the original packet. The L3 layer must be IP. It is optional, but recommended, that the new tunnel encapsulation be parsable by the DPL parser. For more details see the [section on NvTunnelTemplate](#).

### Signature:

```
void nv_set_l3tunnel_underlay<HEADERS, NV_TUNNEL_TEMPLATE, VALUES>(
  inout HEADERS headers, NV_TUNNEL_TEMPLATE tunnel, in VALUES values)
```

### Parameters:

- `headers` [in/out]: The packet headers that are to be encapsulated
- `tunnel` [in]: The `NvTunnelTemplate` object that specifies the underlay headers and the header values.

- `values` [in]: A list expression whose length matches the number of header fields marked "variable" in the `NvTunnelTemplate` object's annotation. The order of the header fields determines which value a field corresponds to, i.e. 3rd header field marked "variable" will be set to the 3rd value in the list. Each value is required to be runtime constant, meaning it may either be an integer value or an action parameter value.

## nv\_mirror

### Description:

Extern function that duplicates the packet, sending each to a different port. Hence this extern is a terminal action.

### Signature:

```
void nv_mirror(in nv_logical_port_t vport, in nv_logical_port_t mirror_port)
```

### Parameters:

- `vport` [in]: Destination for the first copy of the packet
- `mirror_port` [in]: Destination for the second copy of the packet

## nv\_mirror\_vxlan\_v4\_to\_remote

### Description:

Extern function that duplicates the packet, sending each to a different port, with the copy being sent to the mirror port encapsulated with a VXLAN IPV4 underlay

### Signature:

```
void nv_mirror_vxlan_v4_to_remote(in nv_logical_port_t vport, in nv_logical_port_t mirror_port, in bool has_vlan, in nv_mac_addr_t dst_mac, in nv_mac_addr_t src_mac, in nv_vlan_id_t vid, in nv_ipv4_addr_t sip, in nv_ipv4_addr_t dip, in nv_vxlan_id_t vni)
```

## Parameters:

- `vport` [in]: Destination for the unmodified copy of the packet
- `mirror_port` [in]: The destination for the encapsulated copy of the packet.
- `has_vlan` [in]: Specifies whether a VLAN tag is present.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set).
- `sip` [in]: Source IPv4 address.
- `dip` [in]: Destination IPv4 address.
- `vni` [in]: VXLAN Network Identifier.

## nv\_mirror\_vxlan\_v6\_to\_remote

### Description:

Extern function that duplicates the packet, sending each to a different port, with the copy being sent to the mirror port encapsulated with a VXLAN IPv6 underlay.

### Signature:

```
void nv_mirror_vxlan_v6_to_remote(in nv_logical_port_t vport, in
nv_logical_port_t mirror_port, in bool has_vlan, in nv_mac_addr_t
dst_mac, in nv_mac_addr_t src_mac, in nv_vlan_id_t vid, in
nv_ipv6_addr_t sip, in nv_ipv6_addr_t dip, in nv_vxlan_id_t vni)
```

### Parameters:

- `vport` [in]: Destination for the unmodified copy of the packet.
- `mirror_port` [in]: Destination for the encapsulated copy of the packet.



- `has_vlan` [in]: Specifies whether a VLAN tag is present.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set).
- `sip` [in]: Source IPv6 address.
- `dip` [in]: Destination IPv6 address.
- `vni` [in]: VXLAN Network Identifier.

## nv\_mirror\_gre\_v4\_to\_remote

### Description:

Extern function that duplicates the packet, sending each to a different port, with the copy being sent to the mirror port encapsulated with a GRE IPv4 underlay, including optional VLAN tagging and GRE key.

### Signature:

```
void nv_mirror_gre_v4_to_remote(in nv_logical_port_t vport, in
nv_logical_port_t mirror_port, in bool has_vlan, in bool has_key, in
nv_mac_addr_t src_mac, in nv_mac_addr_t dst_mac, in nv_vlan_id_t
vid, in nv_ipv4_addr_t src_ip, in nv_ipv4_addr_t dst_ip, in
NvInnerProtocolType proto_type, in bit<32> key)
```

### Parameters:

- `vport` [in]: Destination for the unmodified copy of the packet.
- `mirror_port` [in]: Destination for the encapsulated copy of the packet.
- `has_vlan` [in]: Specifies if a VLAN tag is present.
- `has_key` [in]: Specifies whether the GRE key is present.

- `src_mac` [in]: Source MAC address.
- `dst_mac` [in]: Destination MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set).
- `src_ip` [in]: Source IPv4 address.
- `dst_ip` [in]: Destination IPv4 address.
- `proto_type` [in]: Protocol type for the inner payload.
- `key` [in]: GRE key (only valid if `has_key` is set).

## nv\_mirror\_gre\_v6\_to\_remote

### Description:

Extern function that duplicates the packet, sending each to a different port, with the copy being sent to the mirror port encapsulated with a GRE IPv6 underlay, including optional VLAN tagging and GRE key.

### Signature:

```
void nv_mirror_gre_v6_to_remote(in nv_logical_port_t vport, in
nv_logical_port_t mirror_port, in bool has_vlan, in bool has_key, in
nv_mac_addr_t src_mac, in nv_mac_addr_t dst_mac, in nv_vlan_id_t
vid, in nv_ipv6_addr_t sip, in nv_ipv6_addr_t dip, in
NvInnerProtocolType proto_type, in bit<32> key)
```

### Parameters:

- `vport` [in]: Destination for the unmodified copy of the packet.
- `mirror_port` [in]: Destination for the encapsulated copy of the packet.
- `has_vlan` [in]: Indicates if a VLAN tag is present.
- `has_key` [in]: Specifies whether the GRE key is utilized.

- `src_mac` [in]: Source MAC address.
- `dst_mac` [in]: Destination MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set).
- `sip` [in]: Source IPv6 address.
- `dip` [in]: Destination IPv6 address.
- `proto_type` [in]: Protocol type for the inner payload.
- `key` [in]: GRE key (only valid if `has_key` is set).

## nv\_mirror\_geneve\_v4\_to\_remote

### Description:

Extern function that duplicates the packet, sending each to a different port, with the copy being sent to the mirror port encapsulated with a Geneve IPv4 underlay, including optional VLAN tagging and Geneve options.

### Signature:

```
void nv_mirror_geneve_v4_to_remote(in nv_logical_port_t vport, in
nv_logical_port_t mirror_port, in bool has_vlan, in bool has_option,
in nv_mac_addr_t dst_mac, in nv_mac_addr_t src_mac, in nv_vlan_id_t
vid, in nv_ipv4_addr_t src_ip, in nv_ipv4_addr_t dst_ip, in bit<6>
opt_fields_len, in bit<1> oam, in bit<1> critical, in
NvInnerProtocolType proto_type, in bit<24> vni, in bit<16>
opt_class, in bit<8> opt_type, in bit<5> opt_len, in bit<32>
opt_data)
```

### Parameters:

- `vport` [in]: Destination for the unmodified copy of the packet.
- `mirror_port` [in]: Destination for the encapsulated copy of the packet.

- `has_vlan` [in]: Indicates if a VLAN tag is present.
- `has_option` [in]: Specifies whether Geneve options are utilized.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set).
- `src_ip` [in]: Source IPv4 address.
- `dst_ip` [in]: Destination IPv4 address.
- `opt_fields_len` [in]: Length of the options field in 4-byte units (must be set to 0 if `has_option` is false, otherwise must be set to 2)
- `oam` [in]: OAM flag.
- `critical` [in]: Critical options flag.
- `proto_type` [in]: Protocol type for the inner payload.
- `vni` [in]: Virtual Network Identifier.
- `opt_class` [in]: Option class (ignored if `has_option` is false)
- `opt_type` [in]: Option type. (ignored if `has_option` is false)
- `opt_len` [in]: Option length in 4-byte units. (ignored if `has_option` is false, otherwise must be 1)
- `opt_data` [in]: Option data. (ignored if `has_option` is false)

## **nv\_mirror\_geneve\_v6\_to\_remote**

### **Description:**

Extern function that duplicates the packet, sending each to a different port, with the copy being sent to the mirror port encapsulated with a Geneve IPv6 underlay, including optional VLAN tagging and Geneve options.

### Signature:

```
void nv_mirror_geneve_v6_to_remote(in nv_logical_port_t vport, in
nv_logical_port_t mirror_port, in bool has_vlan, in bool has_option,
in nv_mac_addr_t dst_mac, in nv_mac_addr_t src_mac, in nv_vlan_id_t
vid, in nv_ipv6_addr_t src_ip, in nv_ipv6_addr_t dst_ip, in bit<6>
opt_fields_len, in bit<1> oam, in bit<1> critical, in
NvInnerProtocolType proto_type, in bit<24> vni, in bit<16>
opt_class, in bit<8> opt_type, in bit<5> opt_len, in bit<32>
opt_data)
```

### Parameters:

- `vport` [in]: Destination for the unmodified copy of the packet.
- `mirror_port` [in]: Destination for the encapsulated copy of the packet.
- `has_vlan` [in]: Indicates if a VLAN tag is present.
- `has_option` [in]: Specifies whether Geneve options are utilized.
- `dst_mac` [in]: Destination MAC address.
- `src_mac` [in]: Source MAC address.
- `vid` [in]: VLAN ID (only valid if `has_vlan` is set).
- `src_ip` [in]: Source IPv6 address.
- `dst_ip` [in]: Destination IPv6 address.
- `opt_fields_len` [in]: Length of the options field in 4-byte units (must be set to 0 if `has_option` is false, otherwise must be set to 2)
- `oam` [in]: OAM flag.
- `critical` [in]: Critical options flag.

- `proto_type` [in]: Protocol type for the inner payload.
- `vni` [in]: Virtual Network Identifier.
- `opt_class` [in]: Option class (ignored if `has_option` is false)
- `opt_type` [in]: Option type. (ignored if `has_option` is false)
- `opt_len` [in]: Option length in 4-byte units. (ignored if `has_option` is false, otherwise must be 1)
- `opt_data` [in]: Option data. (ignored if `has_option` is false)

## nv\_send\_debug\_pkt

### Description:

Extern function that duplicates the original packet, adds pipeline metadata and packet state, then sends it to the DPL Debugger. The original packet continues on the pipeline.

### Signature:

```
void nv_send_debug_pkt(in nv_debug_cookie_t cookie = 8w0x0)
```

### Parameters:

- `cookie` [in]: An optional cookie value can be provided by the user, that will show up in the DPL Debugger with the debug packet

### Note

You must enable debug mode by compiling your DPL program with the `-g` flag in addition calling `nv_send_debug_pkt`. Debug packets can only be processed when the packet direction is from UPLINK towards VF (Rx direction) and only from the primary wire port P0. Placing this extern call in a Tx packet path will result in no debug packet. Debugging packets from the second wire port P1 is currently not supported.

See [DOCA Pipeline Language Developer Tool](#) for more details.

## Tables

### Match Kinds

The DPL compiler supports following match types with the restrictions below:

- Exact match
  - Support for at least 9 keys of DWORD width and 8 keys of byte width. Note that the actual limit is dependent on the HW and may fail only at load time.
  - See [Exact tables in the PSA spec](#) for general guidance and usage from a programmer perspective.
- LPM match
  - Not supported for bool-type keys.
  - Not supported for standard metadata fields.
  - See [LPM tables in the PSA spec](#) for general guidance and usage from a controller perspective.
- Ternary match
  - Not supported for bool-type keys.
  - Not supported for standard metadata fields.
  - See [Ternary tables in the PSA spec](#) for general guidance and usage from a programmer perspective.
- Range match
  - Support up to a maximum of 4 range match keys per P4 table .
  - Support for slices of header fields and std\_meta fields as keys.

- Not supported for slices of user-declared variables and user metadata fields as keys.
- Not supported for keys larger than 32 bits (unless sliced down to 32 bits or less, with some key-specific slice alignment restrictions).
- Range match keys in the same P4 Table with Ternary match keys or LPM keys are currently not supported.
- When using multiple range keys in a single P4 table, some combinations of keys may not be supported. Note that this may fail only at load time.
- See [Range tables in the PSA spec](#) for general guidance and usage from a programmer perspective.

## DPL Key Support

IsValid method call, local variables excluding packet-in struct instance, fixed or flex header fields, and standard metadata can be used as a table key.

However, there are following restrictions:

- A key cannot be used multiple times within one P4 table; this will actively fail during compilation. When IsValid method call is used as a table key, the compiler allocates a register for the key. Thus, using multiple IsValid method calls as keys within a table causes the undefined behavior.
- Complex expressions must use name annotation.

The following table lists the match key support for fixed header fields.

Fixed Header Fields	Bit Width	Exact	LPM	Ternary	Range	Notes
<code>headers.ethernet.dst_addr</code>	48	✓	✓	✓	✓ <u>1</u>	
<code>headers.ethernet.src_addr</code>	48	✓	✓	✓	✓ <u>1</u>	



Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
headers.ethernet.ether_type	16	☐	☐	☐	☐	Last extracted I2 etherType can be matched through std_meta.last_I2_ether_type
headers.vlan.vlan_pcp	3	✓	✓	✓	✓	
headers.vlan.vlan_dei	1	✓	✓	✓	✓	
headers.vlan.vlan_id	12	✓	✓	✓	✓	
headers.vlan.vlan_ether_type	16	☐	☐	☐	☐	Last extracted I2 etherType can be matched through std_meta.last_I2_ether_type
headers.inner_ethernet.dst_addr	48	✓	✓	✓	✓ <u>1</u>	
headers.inner_ethernet.src_addr	48	✓	✓	✓	✓ <u>1</u>	
headers.inner_ethernet.ether_type	16	☐	☐	☐	☐	Last extracted inner I2 etherType can be matched through std_meta.inner_last_I2_ether_type
headers.inner_vlan.vlan_pcp	3	✓	✓	✓	✓	
headers.inner_vlan.vlan_dei	1	✓	✓	✓	✓	
headers.inner_vlan.vlan_id	12	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
<code>headers.inner_vlan.vlan_ether_type</code>	16	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Last extracted inner I2 etherType can be matched through <code>std_meta.inner_last_I2_ether_type</code>
<code>headers.ipv4.version</code> Alias with <code>headers.ipv6.version</code>	4	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.ipv4.ihl</code>	4	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<code>headers.ipv4.diff_serv</code> Alias with <code>headers.ipv6.diff_serv</code>	6	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<code>headers.ipv4.ecn</code> Alias with <code>headers.ipv6.ecn</code>	2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<code>headers.ipv4.total_len</code>	16	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<code>headers.ipv4.identification</code>	16	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<code>headers.ipv4.flags</code>	3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<code>headers.ipv4.frag_offset</code>	13	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
<code>headers.ipv4.ttl</code> Alias with <code>headers.ipv6.hop_limit</code>	8	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
headers.ipv4.protocol Alias with headers.ipv6.next_header	8	✓	✓	✓	✓	
headers.ipv4.hdr_checksum	16	✓	✓	✓	✓	
headers.ipv4.src_addr	32	✓	✓	✓	✓	
headers.ipv4.dst_addr	32	✓	✓	✓	✓	
headers.inner_ipv4.version Alias with headers.inner_ipv6.version	4	☐	☐	☐	☐	
headers.inner_ipv4.ihl	4	✓	✓	✓	✓	
headers.inner_ipv4.diffserv Alias with headers.inner_ipv6.diffserv	6	✓	✓	✓	✓	
headers.inner_ipv4.ecn Alias with headers.inner_ipv6.ecn	2	✓	✓	✓	✓	
headers.inner_ipv4.total_len	16	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
headers.inner_ipv4.identification	16	✓	✓	✓	✓	
headers.inner_ipv4.flags	3	✓	✓	✓	✓	
headers.inner_ipv4.frag_offset	13	✓	✓	✓	✓	
headers.inner_ipv4.ttl Alias with headers.inner_ipv6.hop_limit	8	✓	✓	✓	✓	
headers.inner_ipv4.protocol Alias with headers.inner_ipv6.protocol	8	✓	✓	✓	✓	
headers.inner_ipv4.hdr_checksum	16	✓	✓	✓	✓	
headers.inner_ipv4.src_addr	32	✓	✓	✓	✓	
headers.inner_ipv4.dst_addr	32	✓	✓	✓	✓	
headers.ipv6.flow_label	20	✓	✓	✓	✓	
headers.ipv6.payload_length	16	✓	✓	✓	✓	
headers.ipv6.src_addr	128	✓	✓	✓	✓ <u>2</u>	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
headers.ipv6.dst_addr	128	✓	✓	✓	✓ <u>2</u>	
headers.inner_ipv6.flow_label	20	✓	✓	✓	✓	
headers.inner_ipv6.payload_length	16	✓	✓	✓	✓	
headers.inner_ipv6.src_addr	128	✓	✓	✓	✓ <u>2</u>	
headers.inner_ipv6.dst_addr	128	✓	✓	✓	✓ <u>2</u>	
headers.mpls.label	20	✓	✓	✓	✓	
headers.mpls.tc	3	✓	✓	✓	✓	
headers.mpls.bos	1	✓	✓	✓	✓	
headers.mpls.ttl	8	✓	✓	✓	✓	
headers.inner_mpls.label	20	✓	✓	✓	✓	
headers.inner_mpls.tc	3	✓	✓	✓	✓	
headers.inner_mpls.bos	1	✓	✓	✓	✓	
headers.inner_mpls.ttl	8	✓	✓	✓	✓	
headers.icmp.type	8	✓	✓	✓	✓	
headers.icmp.code	8	✓	✓	✓	✓	
headers.icmp.checksum	16	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
headers.icmp.identifier	16	✓	✓	✓	✓	
headers.icmp.sequence_number	16	✓	✓	✓	✓	
headers.inner_icmp.type	8	✓	✓	✓	✓	
headers.inner_icmp.code	8	✓	✓	✓	✓	
headers.inner_icmp.checksum	16	✓	✓	✓	✓	
headers.inner_icmp.identifier	16	✓	✓	✓	✓	
headers.inner_icmp.sequence_number	16	✓	✓	✓	✓	
headers.icmpv6.type	8	✓	✓	✓	✓	
headers.icmpv6.code	8	✓	✓	✓	✓	
headers.icmpv6.checksum	16	✓	✓	✓	✓	
headers.icmpv6.payload_1	32	✓	✓	✓	✓	
headers.icmpv6.payload_2	32	✓	✓	✓	✓	
headers.inner_icmpv6.type	8	✓	✓	✓	✓	
headers.inner_icmpv6.code	8	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
headers.inner_icmpv6.checksum	16	✓	✓	✓	✓	
headers.inner_icmpv6.payload_1	32	✓	✓	✓	✓	
headers.inner_icmpv6.payload_2	32	✓	✓	✓	✓	
headers.tcp.src_port Alias with headers.udp.src_port	16	✓	✓	✓	✓	
headers.tcp.dst_port Alias with headers.udp.dst_port	16	✓	✓	✓	✓	
headers.tcp.seq_no	32	✓	✓	✓	✓	
headers.tcp.ack_no	32	✓	✓	✓	✓	
headers.tcp.data_offset	4	✓	✓	✓	✓	
headers.tcp.res	3	☐	☐	☐	☐	
headers.tcp.nonce_sum	1	✓	✓	✓	✓	
headers.tcp.ecn	2	✓	✓	✓	✓	
headers.tcp.flags	6	✓	✓	✓	✓	
headers.tcp.window	16	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
<code>headers.tcp.checksum</code>	16	✓	✓	✓	✓	
<code>headers.tcp.urgent_ptr</code>	16	✓	✓	✓	✓	
<code>headers.inner_tcp.src_port</code> Alias with <code>headers.inner_udp.src_port</code>	16	✓	✓	✓	✓	
<code>headers.inner_tcp.dst_port</code> Alias with <code>headers.inner_udp.dst_port</code>	16	✓	✓	✓	✓	
<code>headers.inner_tcp.seq_no</code>	32	☐	☐	☐	☐	
<code>headers.inner_tcp.ack_no</code>	32	☐	☐	☐	☐	
<code>headers.inner_tcp.data_offset</code>	4	✓	✓	✓	✓	
<code>headers.inner_tcp.res</code>	3	☐	☐	☐	☐	
<code>headers.inner_tcp.nonce_sum</code>	1	✓	✓	✓	✓	
<code>headers.inner_tcp.ecn</code>	2	✓	✓	✓	✓	
<code>headers.inner_tcp.flags</code>	6	✓	✓	✓	✓	



Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
headers.inner_tcp.window	16	☐	☐	☐	☐	
headers.inner_tcp.checksum	16	✓	✓	✓	✓	
headers.inner_tcp.urgent_ptr	16	☐	☐	☐	☐	
headers.udp.length	16	✓	✓	✓	✓	
headers.udp.checksum	16	✓	✓	✓	✓	
headers.inner_udp.length	16	✓	✓	✓	✓	
headers.inner_udp.checksum	16	✓	✓	✓	✓	
headers.gre.checksum_present	1	✓	✓	✓	✓	
headers.gre.reserved1	1	✓	✓	✓	✓	
headers.gre.key_present	1	✓	✓	✓	✓	
headers.gre.sequence_present	1	✓	✓	✓	✓	
headers.gre.reserved2	4	✓	✓	✓	✓	
headers.gre.reserved3	5	✓	✓	✓	✓	
headers.gre.version	3	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
<code>headers.gre.protocol</code>	16	✓	✓	✓	✓	
<code>headers.nvgre_vsid.vsid</code>	24	✓	✓	✓	✓	
<code>headers.nvgre_vsid.flow_id</code>	8	✓	✓	✓	✓	
<code>headers.esp.security_parameters.index</code>	32	✓	✓	✓	✓	
<code>headers.esp.sequence_number</code>	32	✓	✓	✓	✓	
<code>headers.esp.next_header</code>	8	☐	☐	☐	☐	
<code>headers.psp.next_header</code>	8	✓	✓	✓	✓	
<code>headers.psp.hdr_ext_len</code>	8	✓	✓	✓	✓	
<code>headers.psp.crypt_offset</code>	8	✓	✓	✓	✓	
<code>headers.psp.needs_sampling</code>	1	✓	✓	✓	✓	
<code>headers.psp.drop</code>	1	✓	✓	✓	✓	
<code>headers.psp.version</code>	1	✓	✓	✓	✓	
<code>headers.psp.has_virtualization_key</code>	1	✓	✓	✓	✓	
<code>headers.psp.one_1</code>	1	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LPM	Ternary	Range	Notes
headers.psp.security_parameters_index	32	✓	✓	✓	✓	
headers.psp.initialization_vector	64	✓	✓	✓	✓	
headers.psp.virtualization_key_high	32	✓	✓	✓	✓	
headers.psp.virtualization_key_lower	32	✓	✓	✓	✓	
headers.vxlan.reserved1	4	✓	✓	✓	✓	
headers.vxlan.vni_valid	1	✓	✓	✓	✓	
headers.vxlan.reserved2	3	✓	✓	✓	✓	
headers.vxlan.reserved3	24	✓	✓	✓	✓	
headers.vxlan.vni	24	✓	✓	✓	✓	
headers.vxlan.reserved4	8	✓	✓	✓	✓	
headers.vxlan_gpe.reserved1	4	✓	✓	✓	✓	
headers.vxlan_gpe.vni_valid	1	✓	✓	✓	✓	
headers.vxlan_gpe.reserved2	3	✓	✓	✓	✓	

Fixed Header Fields	Bit Width	Exact	LP	Ternary	Range	Notes
headers.vxlan_gpe.reserved3	16	✓	✓	✓	✓	
headers.vxlan_gpe.next_proto	8	✓	✓	✓	✓	
headers.vxlan_gpe.vni	24	✓	✓	✓	✓	
headers.vxlan_gpe.reserved4	8	✓	✓	✓	✓	
headers.geneve.ver	2	✓	✓	✓	✓	
headers.geneve.opt_len	6	✓	✓	✓	✓	
headers.geneve.o	1	✓	✓	✓	✓	
headers.geneve.c	1	✓	✓	✓	✓	
headers.geneve.reserved1	6	✓	✓	✓	✓	
headers.geneve.protocol_type	16	✓	✓	✓	✓	
headers.geneve.vni	24	✓	✓	✓	✓	
headers.geneve.reserved2	8	✓	✓	✓	✓	

1. Must be a slice of upper 16 bits or lower 32 bits                    

2. Must be a slice of one of 4 dwords                    

The following table lists the support for "isValid()" for fixed header.

Header IsValid	Notes
ethernet	
inner_ethernet	
vlan	VLAN tag is CVLAN
inner_vlan	Inner VLAN tag is CVLAN
svlan	VLAN tag is SVLAN
inner_svlan	Inner VLAN tag is SVLAN
ipv4	L3 header protocol is IPv4
inner_ipv4	Inner L3 header protocol is IPv4
ipv6	L3 header protocol is IPv6
inner_ipv6	Inner L3 header protocol is IPv6
mpls	MPLS layer 0 is valid after passing all HW checks
inner_mpls	Inner MPLS layer 0 is valid after passing all HW checks
mpls1	MPLS layer 1 is valid after passing all HW checks
inner_mpls1	Inner MPLS layer 1 is valid after passing all HW checks
mpls2	MPLS layer 2 is valid after passing all HW checks
inner_mpls2	Inner MPLS layer 2 is valid after passing all HW checks
mpls3	MPLS layer 3 is valid after passing all HW checks
inner_mpls3	Inner MPLS layer 3 is valid after passing all HW checks
mpls_4	MPLS layer 4 is valid after passing all HW checks
inner_mpls4	Inner MPLS layer 4 is valid after passing all HW checks
tcp	L4 type is TCP
inner_tcp	Inner L4 type is TCP
udp	L4 type is UDP
inner_udp	Inner L4 type is UDP

Header IsValid	Notes
gre	<ul style="list-style-type: none"> <li>L3 layer is valid after passing all HW checks</li> <li>IPv4 protocol field value is 47</li> </ul>
nvgre	<ul style="list-style-type: none"> <li>L3 is OK</li> <li>IPv4 protocol is 47</li> <li>GRE key is present</li> <li>GRE protocol type is 25944</li> </ul>
esp	L4 type is IPSEC
icmp	L4 type is ICMP
inner_icmp	Inner I4 type is ICMP
icmpv6	L4 type is ICMP
inner_icmpv6	Inner I4 type is ICMP
vxlan	<ul style="list-style-type: none"> <li>L4 type is UDP</li> <li>headers.udp.dst_port is 4789</li> </ul>
vxlan_gpe	<ul style="list-style-type: none"> <li>L4 type is UDP</li> <li>headers.udp.dst_port is 4790</li> </ul>
geneve	<ul style="list-style-type: none"> <li>L4 type is UDP</li> <li>headers.udp.dst_port is 6081</li> </ul>

The following table lists the match key support for BlueField standard metadata fields.

Standard Metadata Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
ingress_port	32	✓	☐	☐	☐	The P4 port ID that the packet ingressed the pipeline
eth_to_fcs_packet_len	16	✓	☐	☐	✓	Length of the packet from L2 start to FCS, in bytes
is_l2_ok	1	✓	☐	☐	✓	L2 layer is valid after passing all HW checks

Standard Metadata Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
l2_type	2	✓	☐	☐	✓	<ul style="list-style-type: none"> <li>0: L2_TYPE_UNICAST</li> <li>1: L2_TYPE_MULTICAST</li> <li>2: L2_TYPE_BROADCAST</li> </ul>
last_l2_ether_type	16	✓	✓	✓	✓	Last extracted value of etherType within ethernet header or VLAN tags
vlan_type	2	✓	☐	☐	✓	<ul style="list-style-type: none"> <li>0: VLAN_TYPE_NONE</li> <li>1: VLAN_TYPE_SVLAN</li> <li>2: VLAN_TYPE_CVLAN</li> </ul>
is_l3_ok	1	✓	☐	☐	✓	L3 layer is valid after passing all hardware checks
l3_type	2	✓	☐	☐	✓	<ul style="list-style-type: none"> <li>0: L3_TYPE_NONE</li> <li>1: L3_TYPE_IPV4</li> <li>2: L3_TYPE_IPV6</li> </ul>
is_ip_fragmented	1	✓	☐	☐	✓	
is_ipv4_checksum_ok	1	✓	☐	☐	✓	IPv4 layer checksum is valid
is_l4_ok	1	✓	☐	☐	✓	L4 layer (TCP/UDP) is valid after passing all hardware checks
l4_type	2	✓	☐	☐	✓	<ul style="list-style-type: none"> <li>0: L4_TYPE_NONE</li> <li>1: L4_TYPE_TCP</li> <li>2: L4_TYPE_UDP</li> <li>3: L4_TYPE_IPSEC</li> </ul>
l4_type_ext	4	✓	☐	☐	✓	<ul style="list-style-type: none"> <li>0: L4_TYPE_EXT_NONE</li> <li>1: L4_TYPE_EXT_TCP</li> <li>2: L4_TYPE_EXT_UDP</li> <li>3: L4_TYPE_EXT_IPSEC</li> </ul>
is_l4_checksum_ok	1	✓	☐	☐	✓	L4 layer (TCP/UDP) checksum is valid

Standard Metadata Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
l4_src_port	16	✓	✓	✓	✓	
l4_dst_port	16	✓	✓	✓	✓	
encap_type	2	✓	□	□	✓	<ul style="list-style-type: none"> <li>0: ENCAP_TYPE_NONE</li> <li>1: ENCAP_TYPE_L2_TUNNEL</li> <li>2: ENCAP_TYPE_L3_TUNNEL</li> <li>3: ENCAP_TYPE_ROCE</li> </ul>
ipsec_layer	2	✓	□	□	✓	<ul style="list-style-type: none"> <li>0: IPSEC_TYPE_NONE</li> <li>1: IPSEC_TYPE_OVER_IP</li> <li>2: IPSEC_TYPE_OVER_UDP</li> </ul>
ipsec_syndrome	8	✓	□	□	✓	0: CRYPTO_OK
psp_syndrome	8	✓	□	□	✓	0: CRYPTO_OK
is_inner_l2_ok	1	✓	□	□	✓	Inner L2 layer is valid after passing all HW checks
inner_l2_type	2	✓	□	□	✓	<ul style="list-style-type: none"> <li>0: L2_TYPE_UNICAST</li> <li>1: L2_TYPE_MULTICAST</li> <li>2: L2_TYPE_BROADCAST</li> </ul>
inner_last_l2_ether_type	16	✓	✓	✓	✓	Last extracted value of etherType within inner ethernet header and VLAN tags
inner_vlan_type	2	✓	✓	✓	✓	<ul style="list-style-type: none"> <li>0: VLAN_TYPE_NONE</li> <li>1: VLAN_TYPE_SVLAN</li> <li>2: VLAN_TYPE_CVLAN</li> </ul>
is_inner_l3_ok	1	✓	□	□	✓	Inner I3 layer (IPv4/IPv6) is valid after passing all HW checks
inner_l3_type	2	✓	□	□	✓	<ul style="list-style-type: none"> <li>0: L3_TYPE_NONE</li> <li>1: L3_TYPE_IPV4</li> <li>2: L3_TYPE_IPV6</li> </ul>



Standard Metadata Fields	Bit Width	Exact	LP M	Ternary	Range	Notes
is_inner_ipv4_checksum_ok	1	✓	✓	✓	✓	
is_inner_l4_ok	1	✓	☐	☐	✓	Inner L4 layer (TCP/UDP) is valid after passing all HW checks
inner_l4_type	2	✓	☐	☐	✓	<ul style="list-style-type: none"> <li>0: L4_TYPE_NONE</li> <li>1: L4_TYPE_TCP</li> <li>2: L4_TYPE_UDP</li> <li>3: L4_TYPE_IPSEC</li> </ul>
inner_l4_type_ext	4	✓	☐	☐	✓	<ul style="list-style-type: none"> <li>0: L4_TYPE_EXT_NONE</li> <li>1: L4_TYPE_EXT_TCP</li> <li>2: L4_TYPE_EXT_UDP</li> <li>3: L4_TYPE_EXT_IPSEC</li> </ul>
is_inner_l4_checksum_ok	1	✓	☐	☐	✓	
inner_l4_src_port	16	✓	✓	✓	✓	
inner_l4_dst_port	16	✓	✓	✓	✓	
random_value	16	✓	✓	✓	✓	
ut_clock	64	✓	✓	✓	✓	UTC time stamp <ul style="list-style-type: none"> <li>seconds in the upper 32 bits</li> <li>nano-seconds in the lower 32 bits</li> </ul>
fr_clock	64	✓	✓	✓	✓	Free running clock, in units of 1/device frequency
source_qp	24	✓	☐	☐	✓	

## Size

Each P4 table can specify a size attribute. In the DOCA TA, this size represents the maximum number of entries, excluding miss (default) entry. A user-provided size value must be non-zero value. If the size attribute is not present, a default size of 128 is assigned. Note that the memory model of BlueField means that the size parameter is not a guaranteed size, and runtime insertions may fail even if the entry count is below the maximum.

## Default Action

The table can contain default action. When there are no matching entries within this table, it will execute the default action. When the default action is not provided, the missing entry will execute NoAction action (continue to the next logical step in the pipeline).

```
table ipv4_fowarding_table {
    key = {
        headers.ipv4.src_addr : exact;
    }
    actions = {
        forward;
        NoAction;
    }
    default_action = forward(1);
}
```

## Const Entries

Const entries are optional field of P4 table, which can be used to insert entries when loading the blob to BlueField. For LPM or ternary matching, the unmasked bits must be zeros.

```
table ipv4_fowarding_table {
    key = {
```

```

        headers.ipv4.src_addr : lpm;
    }
    actions = {
        forward;
        NoAction;
    }
    default_action = forward(1);
    const entries = {
        (32w0x11111100 &&& 32w0xFFFFFFFF00) : forward(2);
        (32w0x11111100 &&& 36w0xFFFFFFFF000) : forward(3);
    }
}

```

## Packet IO

Controller Packet in (pipeline to controller) and Packet out (controller to pipeline) are supported in the DOCA TA. To send a packet to the controller with metadata:

- Create a struct with the annotation `@nv_controller_metadata("packet_in")`
- The struct must be exactly 32 bits in width (user must insert padding as needed)
- Use the extern action

```
extern void nv_send_to_controller<PACKET_META>(in PACKET_META
pkt_in_meta);
```

To check if a packet that has ingress the pipeline from the controller:

- add as a match key "`std_meta.ingress_port`" to a P4 table key set
- add an entry that specifies the key value to be the P4 Port ID defined in the [DPL Port ID assignment](#) step.

## P4 Language Support in DPL

This section contains information on the P4 language features that are supported, known issues and deviations from the [P4 Language Specification](#). All references to the spec refer to P4<sub>16</sub> language specification version *v1.2.4*.

## Introduction

The NVIDIA BlueField networking platform (DPU or SuperNIC) is a specialized Data Processing Unit engineered to significantly enhance data center performance. It achieves this by efficiently offloading, accelerating, and isolating critical tasks related to networking, storage, and security. BlueField-3 merges the advantages of dedicated accelerators with the versatility of general-purpose processors, all within an ASIC-based system-on-a-chip architecture. It boasts impressive connectivity speeds of up to 400G, making it an ideal choice for environments demanding high levels of AI and high-performance computing capabilities.

To ensure BlueField delivers on its promise of optimized functionality, its P4 language implementation has been tailored to leverage the strengths of high-performance ASIC hardware. This strategic focus on performance and efficiency, however, means that BlueField supports DOCA Pipeline Language (DPL), which does not include every feature outlined in the P4 language specification. This is primarily due to the inherent differences in flexibility and programmability between ASICs and other types of hardware, such as CPUs and FPGAs, and the unique pipeline model of the BlueField DPU. This section outlines the P4 Language features that are currently supported in this release of the NVIDIA DPL compiler.

### **Note**

This document refers to features as being **supported** or **unsupported**. Supported features have been tested and should work according to the P4<sub>16</sub> language specification, subject to any caveats described in this document. Unsupported features have not been fully tested and should not be relied upon. In most cases, the compiler will reject programs that use unsupported features. However, in some cases the compiler may accept a program that uses unsupported features if the feature is not necessary to implement the program. For example, if a program contains an expression that uses an unsupported operator but its operands can be computed at compile-time, the compiler may choose to compute the value of the expression at compile-time and accept the program. This behavior

should not be used as an indicator of whether the feature is supported.

## P4 Language Features

### Identifiers

Identifiers starting with `__` are reserved for internal compiler use. Otherwise, identifiers described in the P4 language spec section *6.4.1. Identifiers* are allowed.

### Data Types

See *Operators* section for support of operations on values with these types.

- Bool
  - Supported
- Arbitrary-precision Integer
  - Supported only for literals. See spec [7.1.6.5. Arbitrary-precision integers](#).
- Signed integer
  - Unsupported
- Strings literals
  - No operations are allowed or validity checks performed. See spec [6.4.3.3](#).
- Bit strings
  - Supported, limited by available hardware resources

### Derived Types

- Enum

- Enumeration types are supported as described in section [7.2.1](#) of the P4<sub>16</sub> spec, allowing the P4 programmer to either specify an underlying representation or allow the compiler to choose the representation. Note that the set of allowed types for the underlying representation is limited to those otherwise supported by the DPL compiler.
- Header
  - Header types are supported as described in section [7.2.2](#) of the P4<sub>16</sub> spec using field types otherwise supported by the DPL compiler with the exception of `varbit<>` fields. Variable-length header types are supported only using the `NvOptionParser` extern type.
- Header stacks
  - Unsupported
- Structs
  - Struct types are supported as described in section [7.2.5](#) of the P4<sub>16</sub> spec using field types otherwise supported by the DPL compiler.
- Unions
  - Unsupported
- Tuple/List
  - Tuple types are supported as described in section [7.2.6](#) of the P4<sub>16</sub> spec using component types otherwise supported by the DPL compiler.
- Extern types
  - Extern types, including both extern functions and extern objects, as described in section [7.2.9](#) of the P4<sub>16</sub> spec are supported only for those declared in the P4 headers distributed with the DPL compiler. P4 programs cannot declare additional extern types.
- Type specialization
  - Type specialization is supported as described in section [7.2.10](#) of the P4<sub>16</sub> spec.

# Statements and Expressions

- Assignment
  - An L-value cannot be used in a method call expression, packet out metadata, flex-header field, standard metadata field, or as an action parameter. Not all fixed header fields can be an L-value of an assignment statement. Please refer to the chart below.
- Conditional
  - Conditional statements are only supported within control apply blocks. Its expression must evaluate to a bit or bool type.
- switch statement
  - The switch statement is only supported within control apply block. Its expression must evaluate to a bool type.
  - The compiler supports empty switch statement, fall through, default case, and non-default cases. See spec section [11.7 Switch statement](#) for details.

The following tables describe the compiler support for expressions using the built in header fields and standard metadata as L-values and R-values. Note, this is separate of header fields that can be used as match keys.

## Note

In the default hardware parser, some fields that are mutually exclusive are extracted to the same buffer location (referred to in the table as an alias). Assignments to and copy from these fields can use either of the aliased field names.

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.ethernet.dst_addr</code>	✓	✓	

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.ethernet.src_addr</code>	✓	✓	
<code>headers.ethernet.ether_type</code>	☐	☐	Last extracted outer <code>etherType</code> value
<code>headers.vlan.vlan_pcp</code>	☐	☐	
<code>headers.vlan.vlan_dei</code>	☐	☐	
<code>headers.vlan.vlan_id</code>	✓	✓	
<code>headers.vlan.vlan_ether_type</code>	☐	☐	Last extracted outer <code>etherType</code> value
<code>headers.inner_ethernet.dst_addr</code>	☐	✓	
<code>headers.inner_ethernet.src_addr</code>	☐	✓	
<code>headers.inner_ethernet.ether_type</code>	☐	☐	Last extracted inner <code>etherType</code> value
<code>headers.inner_vlan.vlan_pcp</code>	☐	☐	
<code>headers.inner_vlan.vlan_dei</code>	☐	☐	
<code>headers.inner_vlan.vlan_id</code>	☐	☐	
<code>headers.inner_vlan.vlan_ether_type</code>	☐	☐	Last extracted inner <code>etherType</code> value
<code>headers.ipv4.version</code> Alias with <code>headers.ipv6.version</code>	☐	☐	
<code>headers.ipv4.ihl</code>	☐	☐	
<code>headers.ipv4.diffserv</code> Alias with <code>headers.ipv6.diffserv</code>	☐	✓	Can be set through <code>nv_set_ip_dscp</code> extern



Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.ipv4.ecn</code> Alias with <code>headers.ipv6.ecn</code>	☐	✓	Can be set through <code>nv_set_ip_ecn</code> extern
<code>headers.ipv4.total_len</code>	☐	☐	Value is write only by hardware
<code>headers.ipv4.identificati on</code>	☐	☐	
<code>headers.ipv4.flags</code>	☐	☐	
<code>headers.ipv4.frag_offset</code>	☐	☐	
<code>headers.ipv4.ttl</code> Alias with <code>headers.ipv6.hop_limit</code>	☐	✓	Can be set through <code>nv_set_ip_ttl</code> extern
<code>headers.ipv4.protocol</code> Alias with <code>headers.ipv6.next_header</code>	☐	☐	
<code>headers.ipv4.hdr_checksum</code>	☐	☐	Value is write only by hardware
<code>headers.ipv4.src_addr</code>	✓	✓	
<code>headers.ipv4.dst_addr</code>	✓	✓	
<code>headers.inner_ipv4.versio n</code> Alias with <code>headers.inner_ipv6.versio n</code>	☐	☐	
<code>headers.inner_ipv4.ihl</code>	☐	☐	
<code>headers.inner_ipv4.diffse rv</code> Alias with <code>headers.inner_ipv6.diffse rv</code>	☐	✓	
<code>headers.inner_ipv4.ecn</code> Alias with <code>headers.inner_ipv6.ecn</code>	☐	☐	

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.inner_ipv4.total_len</code>	☐	☐	
<code>headers.inner_ipv4.identification</code>	☐	☐	
<code>headers.inner_ipv4.flags</code>	☐	☐	
<code>headers.inner_ipv4.frag_offset</code>	☐	☐	
<code>headers.inner_ipv4.ttl</code> Alias with <code>headers.inner_ipv6.hop_limit</code>	☐	✓	
<code>headers.inner_ipv4.protocol</code> Alias with <code>headers.inner_ipv6.protocol</code>	☐	☐	
<code>headers.inner_ipv4.hdr_checksum</code>	☐	☐	
<code>headers.inner_ipv4.src_addr</code>	☐	✓	
<code>headers.inner_ipv4.dst_addr</code>	☐	✓	
<code>headers.ipv6.flow_label</code>	☐	☐	
<code>headers.ipv6.payload_length</code>	☐	☐	
<code>headers.ipv6.src_addr</code>	✓	✓	
<code>headers.ipv6.dst_addr</code>	✓	✓	
<code>headers.inner_ipv6.flow_label</code>	☐	☐	

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.inner_ipv6.payload_length</code>	☐	☐	
<code>headers.inner_ipv6.src_addr</code>	☐	✓	
<code>headers.inner_ipv6.dst_addr</code>	☐	✓	
<code>headers.mpls.label</code>	✓	✓	
<code>headers.mpls.tc</code>	✓	✓	
<code>headers.mpls.bos</code>	✓	✓	
<code>headers.mpls.ttl</code>	✓	✓	
<code>headers.inner_mpls.label</code>	☐	✓	
<code>headers.inner_mpls.tc</code>	☐	✓	
<code>headers.inner_mpls.bos</code>	☐	✓	
<code>headers.inner_mpls.ttl</code>	☐	✓	
<code>headers.icmp.type</code>	☐	☐	
<code>headers.icmp.code</code>	☐	☐	
<code>headers.icmp.checksum</code>	☐	☐	
<code>headers.icmp.identifier</code>	☐	☐	
<code>headers.icmp.sequence_number</code>	☐	☐	
<code>headers.inner_icmp.type</code>	☐	☐	
<code>headers.inner_icmp.code</code>	☐	☐	
<code>headers.inner_icmp.checksum</code>	☐	☐	
<code>headers.inner_icmp.identifier</code>	☐	☐	

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.inner_icmp.sequence_number</code>	☐	☐	
<code>headers.icmpv6.type</code>	☐	☐	
<code>headers.icmpv6.code</code>	☐	☐	
<code>headers.icmpv6.checksum</code>	☐	☐	
<code>headers.icmpv6.payload_1</code>	☐	☐	
<code>headers.icmpv6.payload_2</code>	☐	☐	
<code>headers.inner_icmpv6.type</code>	☐	☐	
<code>headers.inner_icmpv6.code</code>	☐	☐	
<code>headers.inner_icmpv6.checksum</code>	☐	☐	
<code>headers.inner_icmpv6.payload_1</code>	☐	☐	
<code>headers.inner_icmpv6.payload_2</code>	☐	☐	
<code>headers.tcp.src_port</code> Alias with <code>headers.udp.src_port</code>	☐	✓	Can be set through <code>nv_set_l4_src_port</code> extern
<code>headers.tcp.dst_port</code> Alias with <code>headers.udp.dst_port</code>	☐	✓	Can be set through <code>nv_set_l4_dst_port</code> extern
<code>headers.tcp.seq_no</code>	✓	✓	
<code>headers.tcp.ack_no</code>	✓	✓	
<code>headers.tcp.data_offset</code>	☐	☐	
<code>headers.tcp.res</code>	☐	☐	
<code>headers.tcp.nonce_sum</code>	☐	☐	

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.tcp.ecn</code>	✓	✓	
<code>headers.tcp.flags</code>	✓	✓	
<code>headers.tcp.window</code>	☐	☐	
<code>headers.tcp.checksum</code>	☐	☐	
<code>headers.tcp.urgent_ptr</code>	☐	☐	
<code>headers.inner_tcp.src_port</code> Alias with <code>headers.inner_udp.src_port</code>	☐	✓	
<code>headers.inner_tcp.dst_port</code> Alias with <code>headers.inner_udp.dst_port</code>	☐	✓	
<code>headers.inner_tcp.seq_no</code>	☐	✓	
<code>headers.inner_tcp.ack_no</code>	☐	✓	
<code>headers.inner_tcp.data_of fset</code>	☐	☐	
<code>headers.inner_tcp.res</code>	☐	☐	
<code>headers.inner_tcp.nonce_s um</code>	☐	☐	
<code>headers.inner_tcp.ecn</code>	☐	✓	
<code>headers.inner_tcp.flags</code>	☐	✓	
<code>headers.inner_tcp.window</code>	☐	☐	
<code>headers.inner_tcp.checksu m</code>	☐	☐	

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.inner_tcp.urgent_ptr</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.udp.length</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.udp.checksum</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.inner_udp.length</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.inner_udp.checksum</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.checksum_present</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.reserved1</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.key_present</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.sequence_present</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.reserved2</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.reserved3</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.version</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.gre.protocol</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.nvgre_vsid.vsid</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.nvgre_vsid.flow_id</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.esp.security_parameters.index</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.esp.sequence_number</code>	<input type="checkbox"/>	<input type="checkbox"/>	
<code>headers.esp.next_header</code>	<input type="checkbox"/>	<input type="checkbox"/>	Value set by hardware after decryption
<code>headers.psp.next_header</code>	<input type="checkbox"/>	<input type="checkbox"/>	Value set by hardware after

Fixed Header Fields	Assignab le	Copyabl e	Notes
			decryption
headers.psp.hdr_ext_len	☐	☐	
headers.psp.crypt_offset	☐	☐	
headers.psp.needs_samplin g	☐	☐	
headers.psp.drop	☐	☐	
headers.psp.version	☐	☐	
headers.psp.has_virtualiz ation_key	☐	☐	
headers.psp.one_1	☐	☐	
headers.psp.security_para meters_index	☐	☐	
headers.psp.initializatio n_vector	☐	☐	
headers.psp.virtualizatio n_key_high	☐	☐	
headers.psp.virtualizatio n_key_low	☐	☐	
headers.vxlan.reserved1	☐	☐	
headers.vxlan.vni_valid	☐	☐	
headers.vxlan.reserved2	☐	☐	
headers.vxlan.reserved3	☐	☐	
headers.vxlan.vni	✓	✓	
headers.vxlan.reserved4	☐	☐	
headers.vxlan_gpe.reserve d1	☐	☐	

Fixed Header Fields	Assignab le	Copyabl e	Notes
<code>headers.vxlan_gpe.vni_val id</code>	☐	☐	
<code>headers.vxlan_gpe.reserve d2</code>	☐	☐	
<code>headers.vxlan_gpe.reserve d3</code>	☐	☐	
<code>headers.vxlan_gpe.next_pr oto</code>	☐	☐	
<code>headers.vxlan_gpe.vni</code>	✓	✓	
<code>headers.vxlan_gpe.reserve d4</code>	☐	☐	
<code>headers.geneve.ver</code>	☐	☐	
<code>headers.geneve.opt_len</code>	☐	☐	
<code>headers.geneve.o</code>	☐	☐	
<code>headers.geneve.c</code>	☐	☐	
<code>headers.geneve.reserved1</code>	☐	☐	
<code>headers.geneve.protocol_t ype</code>	☐	☐	
<code>headers.geneve.vni</code>	✓	✓	
<code>headers.geneve.reserved2</code>	☐	☐	

All the fields of BlueField standard metadata are read only. The following table outlines the current support for using a standard metadata field as an R-value in an expression.



Standard Metadata Fields	Copyable	Notes
<code>ingress_port</code>	✓	
<code>eth_to_fcs_packet_len</code>	☐	
<code>is_l2_ok</code>	☐	
<code>l2_type</code>	☐	
<code>last_l2_ether_type</code>	✓	Last extracted value of ether type with

Standard Metadata Fields	Copyable	Notes
		i n e t h e r n e t h e a d e r o r V L A N t a g s
vlan_type	☐	
is_l3_ok	☐	
l3_type	☐	
is_ip_fragmented	☐	
is_ipv4_checksum_ok	☐	
is_l4_ok	☐	
l4_type	☐	

Standard Metadata Fields	Copyable	Notes
l4_type_ext	☐	
is_l4_checksum_ok	☐	
l4_src_port	☐	
l4_dst_port	☐	
encap_type	☐	R O C E n o t c u r r e n t l y s u p p o r t e d
ipsec_layer	☐	
ipsec_syndrome	☐	V a l i d o

Standard Metadata Fields	Copyable	Notes
		n l y a f t e r h a r d w a r e e n c r y p t/ d e c r y p t
<div data-bbox="164 1549 407 1591" style="border: 1px solid gray; border-radius: 5px; padding: 2px;">             psp_syndrome           </div>	<input data-bbox="800 1549 824 1581" type="checkbox"/>	V a l i d o n l y a f t

Standard Metadata Fields	Copyable	Notes
		e r h a r d w a r e e n c r y p t/ d e c r y p t
is_inner_l2_ok	<input type="checkbox"/>	
inner_l2_type	<input type="checkbox"/>	
inner_last_l2_ether_type	<input checked="" type="checkbox"/>	L a s t e x t r a c t

Standard Metadata Fields	Copyable	Notes
		e d v a l u e o f e t h e r T y p e w i t h i n n e r e t h e r n e t h e a

Standard Metadata Fields	Copyable	Notes
		der or V L A N t a g s
inner_vlan_type	☐	
is_inner_l3_ok	☐	
inner_l3_type	☐	
is_inner_ipv4_checksum_ok	☐	
is_inner_l4_ok	☐	
inner_l4_type	☐	
inner_l4_type_ext	☐	
random_value	☐	
ut_clock	☐	
fr_clock	☐	
source_qp	☐	

## Operators

The P4-16 language specification lists a wide variety of operations that the language accepts for the supported data types (see [Section 8](#)). The table below lists the operators

that are officially supported by the NVIDIA P4 compiler:

Operator	Compile-time value	P4Runtime value	Runtime value	Spec section
Bool && Bool	✓	✓	✓	8.5
Bool    Bool	✓	✓	✓	8.5
Bool == Bool	✓	✓ <u>1</u>	✓ <u>1</u>	8.5
Bool != Bool	✓	✓ <u>1</u>	✓ <u>1</u>	8.5
Bit<W> == Bit<W>	✓	✓ <u>1</u>	✓ <u>1</u>	8.6
Bit<W> != Bit<W>	✓	✓ <u>1</u>	✓ <u>1</u>	8.6
Bit<W> << integer	✓ <u>2</u>	✓ <u>2</u>	✓ <u>2</u>	8.6
Bit<W> >> integer	✓ <u>2</u>	✓ <u>2</u>	✓ <u>2</u>	8.6
Bit<W>[H:L]	✓ <u>3</u>	✓ <u>3</u>	✓ <u>3</u>	8.6
All explicit casts between supported types	✓	✓	✓	8.11.1
All implicit casts between supported types	✓	✓	✓	8.11.2
Bit<W>..Bit<W>	✓ <u>4</u>	✓ <u>4</u>	□	8.15.4
Assignment to user struct fields	✓	✓	✓	8.16
Assignment to packet-in struct fields	✓	✓	✓	8.16
All operations on header fields	✓	✓	✓ <u>5</u>	8.17
Method calls	✓	✓	□	8.20



Operator	Compile-time value	P4Runtime value	Runtime value	Spec section
Function calls with positional args	✓	✓ <u>6</u>	✓ <u>6</u>	8.20
Extern constructor invocations	✓	☐	☐	8.21
Parser constructor invocations	✓	☐	☐	8.21
Control constructor invocations	✓	☐	☐	8.21
Package constructor invocations	✓	☐	☐	8.21

1. LHS or RHS must be compile-time constant
2. RHS must be compile-time constant. See *spec 8.9.2*
3. H and L are subject to the restrictions described in the *spec 8.6*. Assigning to slices (slices as L-values) is not supported. Additionally, slices as R-values are only supported as P4 table keys.
4. Only valid in P4 Table entries
5. Limited to those fields that can be a copy source
6. Limited on a per extern function basis

## Variables

Variables are supported in accordance with the following spec items:

- Constants (*spec 11.1*)
  - "Compile-time known values" are evaluated on a best-effort basis. It is possible that a compile-time known value may not be recognized by the compiler as such.

- Variables (*spec 11.2*)
- Instantiations (*spec 11.3*)
  - Instantiations with abstract methods (*spec 11.3.1*) are allowed in BlueField Target Architecture
  - Named arguments are not supported

Variables may be declared in any of the locations described in (*spec 11.2*) and follow the scope rules described there.

The compiler will emit errors for uninitialized values. In some cases where a struct is partially initialized, only a warning may be produced. In some cases there may be no error emitted when an uninitialized struct field is accessed. The accessed field will then contain an undefined value.

## Control Apply Block

The following statements are supported in a control's apply block:

- `table.apply()` calls
- `if` statement
- `switch` statement
- extern function and method calls
- assignment statements with the supported operators
- the empty statement
- `return` statements

The `exit` statement is not supported.

All supported expressions are allowed within these statements, where applicable.

## Actions

Actions support the same statements as controls except for the following:

- `table.apply()` calls
- Conditional statements - `if` and `switch`

Actions support the same expressions as controls except for the following:

- Boolean logical operators - `&&`, `||`, ternary operator
- Comparisons (`==`, `!=`, etc.)

# DPL Installation Guide

## Introduction

The DPL Development Container (`dp1_dev`) bundles several tools for compiling, controlling and debugging DPL programs.

Pulling the container image can be done using the following command:

```
docker pull nvcr.io/nvidia/doca/dp1_dev:1.0.1-doca2.10.0-host
```

A set of convenience scripts is provided together with `dp1_dev`, one script per tool.

Each script launches each respective tool, requiring fewer command line options to be specified.

## Installing the launch scripts

The `/install.sh` script provided within the `dp1_dev` container copies the launch scripts to the current working directory.

Below is an example command of how it can be used to copy the launch scripts from within the container to a mounted directory of your choosing (specifically `${PWD}` in this example).

```
docker run --rm -it -v ${PWD}:${PWD} -u $(id -u):$(id -g) -w  
${PWD} nvcr.io/nvidia/doca/dpl_dev:1.0.1-doca2.10.0-host  
/install.sh
```

Make sure the copied scripts have execute permission with:

### Info

```
chmod +x scripts/*.sh
```

The above command copies the following scripts to the local host file system:

- Tool launch scripts
  - `dplp4c.sh`
  - `dpl_admin.sh`
  - `p4runtime_sh.sh`
  - `dpl_nspect.sh`
  - `dpl_debugger.sh`
- Configuration script
  - `scripts_config.sh`
- General purpose script for internal use only
  - `scripts_utils.sh`

For a detailed explanation on each tool, see [DOCA Pipeline Language Developer Tool](#).

## Configuring launch scripts

For your convenience, the `scripts_config.sh` file can be used to avoid specifying commonly used arguments when executing the scripts.

By editing the following file you can avoid passing the `-i <DPL Dev image name:tag>` and `-a <DPL Runtime daemon address:port>` to the various tools.

```
#!/bin/bash

# This configuration file defines various arguments used by the launch scripts.
# These variables are optional, but defining them in advance simplifies the use of the launch scripts.

# Image name and tag that contains the development tools
# Example: nvcr.io/nvidia/doca/dpl_dev:<tag>
DEV_IMAGE=nvcr.io/nvidia/doca/dpl_dev:1.0.1-doca2.10.0-host

# IP where the DPL Runtime daemon (RTD) is running
# Example: 192.168.1.100
# DPL_RTD_IP=

# IP where the DPL RTD is running and the port for the P4 Runtime Server
# Example: ${DPL_RTD_IP}:9559
# P4RT_ADDRESS=${DPL_RTD_IP}:9559

# IP where the DPL RTD is running and the port for the DPL Admin
# Example: ${DPL_RTD_IP}:9600
# DPL_ADMIN_ADDRESS=${DPL_RTD_IP}:9600

# IP where the DPL RTD is running and the port for the DPL Nspect
# Example: ${DPL_RTD_IP}:9560
# DPL_NSPECT_ADDRESS=${DPL_RTD_IP}:9560
```

## Next steps

1. [Compiling DPL Applications](#) that you wish to examine using the developer tools.
2. [Loading DPL Applications](#) can now be done with the compilation outputs.
3. Use the [DOCA Pipeline Language Developer Tool](#) to further examine the DPL programs correct operation.

# Compiling DPL Applications

This section describes how to use the NVIDIA DPL compiler to compile DPL applications.

## Introduction

DOCA Programming Language (DPL) applications are an educational resource provided as a guide on how to program on the NVIDIA® BlueField® networking platform (DPU or SuperNIC) using DPL.

The NVIDIA DPL compiler (dplp4c.sh) is provided as an executable in a self contained docker image. The docker image is available to partners by logging into the [NVIDIA NGC](#) system and downloading the most recent [DPL docker containers](#).

## Installation

Please refer to the [DPL Runtime Service](#) for the system requirements and [DPL Installation Guide](#) on how to install the compiler. For information on the development environment, refer to the [DOCA Pipeline Language Model](#) and the [DPL System Overview](#).

## Dependencies

- BlueField-3 is required
- Ubuntu 22.04 hosts (x86) or greater
- A P4Runtime controller based on API version 1.3.0 or greater

## Compiling Applications

The DPL compiler is provided with all its dependencies in the dpl-dev container. For details on how to obtain and install the docker images see [DPL Installation Guide](#) .

A shell script, `dplp4c.sh`, is provided for a convenient way to execute the DPL compiler from the `dpl-dev` container:

```
./dplp4c.sh [--mount dir]* dplp4c_args*
```

- `--mount` – (optional) directory to be mounted into the container
- `dplp4c_args` – arguments to be passed to the compiler (e.g., `sample.p4`)

For example, compiling the "hello\_packet" DPL example is as simple as:

```
local-user@vm-1:~/p4-samples/dpu/hello_packet$ dplp4c.sh hello_packet.p4
Generating compiler output in "_out"
updating: MANIFEST.json (deflated 40%)
updating: hello_packet.program.json (deflated 87%)
updating: hello_packet.debug.json (deflated 96%)
```

The following files in this example are produced in the `_out` directory:

<code>compiler.log</code>	Log of any compiler warnings or errors to the specified program
<code>hello_packet.dplconfig</code>	Binary blob, containing all the data needed for the DPL Runtime daemon to load the program
<code>hello_packet.p4info.txt</code>	P4Runtime protobuf file, in text format

Additional NVIDIA® BlueField®-specific arguments include:

```
--help           Print this help message
--version        Print compiler version
-I path          Specify include path
                 (passed to preprocessor)
-D arg=value     Define macro (passed to
preprocessor)
-U arg           Undefine macro (passed to
preprocessor)
```

-E	Preprocess only, do not compile (prints program on stdout)
-M	Output `make` dependency rule only (passed to preprocessor)
-MD	Output `make` dependency rule to file as side effect (passed to preprocessor)
-MF file	With -M, specify output file for dependencies (passed to preprocessor)
-MG	with -M, suppress errors for missing headers (passed to preprocessor)
-MP	with -M, add phony target for each dependency (passed to preprocessor)
-MT target	With -M, override target of the output rule (passed to preprocessor)
-MQ target	Like -Mt, override target but quote special characters (passed to preprocessor)
-g	Enable debugging via DPL Nspect
--nocpp	Skip preprocess, assume input file is already preprocessed.
--Wdisable[=diagnostic]	Disable a compiler diagnostic, or disable all warnings if no diagnostic is specified.
--Winfo[=diagnostic]	Report an info message for a compiler diagnostic.
--Wwarn[=diagnostic]	Report a warning for a compiler diagnostic, or treat all info messages as warnings if no diagnostic is specified.
--Werror[=diagnostic]	Report an error for a compiler diagnostic, or treat all warnings as errors if no diagnostic is specified.
--maxErrorCount errorCount	Set the maximum number of errors to display before failing.
--target target	Compile for the specified target device.



<code>--odir out_directory</code>	Write output to out directory
<code>--enable feature[,feature]* --help</code>	Enable a feature, or comma-separated list of features
<code>--disable feature[,feature]* --help</code>	Disable a feature, or comma-separated list of features

See section [Compiling DPL Applications](#) for examples of how to use the `dplp4c.sh` to compile.

**i Info**

The binary `<sample_name>.` is created under `_out/<sample_name>.dplconfig.`

**i Info**

The P4Runtime file is created under `_out/<sample_name>.p4info.txt.`

## Compiling Applications in Debug Mode

To enable packet debugging, you must ensure the following:

- Compile your DPL program with the `-g` flag. This enables parsing for a debug flow.
- Ensure you have at least one `nv_send_debug_pkt()` call in your DPL program. A debug packet will be emitted where `nv_send_debug_pkt()` is placed, and represent the program state at that point.

Please note packet debugging is only available on the RX direction. Debug packets will not be produced on TX. You may optionally include a "cookie", within your debug extern call, i.e.: `nv_send_debug_pkt(8w0x42)`. This cookie may help differentiate between different debug extern calls, and will appear in your debugging output.

## Loading DPL Applications

NVIDIA DPL programs are deployed to the NVIDIA® BlueField® networking platform (DPU or SuperNIC) using the P4Runtime API. This allows platform-independent, standards-based integration with SDN controllers.

### Introduction

The DPL compiler generates a pipeline binary optimized for execution on BlueField. Pipeline loading and control—such as installing the program and populating P4 tables—are handled via the P4Runtime API, an open and well-defined interface.

The P4Runtime server, running on the BlueField device, enables a P4 Controller to:

- Connect over gRPC
- Set the `ForwardingPipelineConfig` (i.e., install the compiled DPL binary and `p4info`)
- Query the device for its current pipeline config and table state
- Maintain runtime P4 tables as defined in the DPL source

This model enables integration with open-source, proprietary, or custom-built controllers in a standardized way.

### Prerequisites

Before loading a DPL application, ensure the following services and components are properly set up:

- [DPL Runtime Service](#) is running and configured on the BlueField (Arm side). See the [Container Deployment](#) page for setup instructions.
- [DPL Development Container](#) and the `p4runtime_sh.sh` launch script are installed on the host. See the [DPL Installation Guide](#) for more details.

## Loading the Application

The a DPL application can be loaded using:

- A custom P4Runtime controller
- The the [NVIDIA-supplied Python controller](#)
- An [open source controller](#)

In the following example we'll be using the P4Runtime controller bundled within the [DPL Development Container](#).

## Using p4runtime\_sh.sh Script

Running the script with no arguments displays the usage information:

```
usage: p4runtime_sh.sh -i <docker image> -p <program_folder> -a
<dpl_rtd_ip:port> [OPTIONS]
```

## Example Command

```
p4runtime_sh.sh -i doca_p4_dev:latest -p /root/hello_packet/_out
-a 192.168.1.100:9559 --device-id 1000
```

Arguments:

- `-i` - The pulled [DPL Development Container](#) image.
- `-p` - Directory that holds the DPL program compilation outputs ([Compiling DPL Applications](#)).
- `-a` - Address of the DPL Runtime daemon and the P4Runtime port as specified in the DPL Runtime [Service Configuration](#).

- `--device-id` – (Optional) ID of the target device

After successful loading, the script launches an interactive Python shell connected to the P4Runtime server. From here, you can inspect and manipulate tables (see [p4runtime\\_sh Usage](#)).

## P4Runtime Optimizations

DPL table entries are added via a P4Runtime controller, which may run remotely or locally on BlueField.

- The standard gRPC-based controller model supports ~50K rule insertions per second
- For use cases requiring high-speed (1M+) rule insertions, NVIDIA is introducing a bulk insertion API extension to the P4Runtime protobuf specification

### Note

This feature is planned for future releases.

## Sample DPL Applications

This section describes NVIDIA DOCA Pipeline Language sample applications.

### Sample Prerequisites

The DPL program can be run compiled on any host machine with the DPL compiler installed. The resulting binary can be loaded on either the host machine or on an NVIDIA® BlueField®-3 DPU target where the DPL Runtime daemon is running. See to the following documents:

- - [DPL Runtime Service](#) for details on how to install BlueField-related software.

- [DPL Release Notes](#) for any issue you may encounter with the installation, compilation, or execution of DPL samples.

## Running the Sample

Compile the DPL sample program, locate the generated .p4info and .dplconfig files from the compiler output directory, and load them to the DPL Runtime daemon using a P4Runtime controller.

## Sample Applications

### Hello Packet

[This sample program](#) demonstrates a basic match/action pipeline, with a simple match and a drop or send to port action.

### GTP Parsing

[This sample program](#) demonstrates how to add a custom protocol to the native BlueField Platform parser and match on those fields.

### GTP Tunnel Encapsulation

[This sample program](#) demonstrates how to add a custom tunnel encapsulation, GTP-U, that was not natively supported on the BlueField platform.

### Geneve TLV Parsing

[This sample program](#) demonstrates how to add a custom protocol that contains TLVs, and to match on the TLV fields.

### VXLAN Tunnel Gateway

[This sample program](#) demonstrates how to create a tunnel gateway between two different VXLAN domains.

## Connection Tracking

[This sample program](#) shows how to program the data plane portion of a basic TCP connection tracking solution. Packets are sent to a controller to manage the table entries.

## Host-based Networking

[This sample program](#) provides a stretched L2 tenant pipeline with VXLAN encapsulation/decapsulation on top of the BlueField Platform.

# Hello Packet Example

This example demonstrates the most basic pipeline in the DOCA target architecture. It consists of:

- Match on L2 source MAC address
- Action to either forward to a P4 port ID, drop, or no action

## Sample Code

The following `#includes` are required for every DPL program. They define the DOCA target architecture. In particular, they contain things like the default DOCA parser, default headers struct, etc. Note that the symbol names are prefixed with "nv\_" and are reserved for NVIDIA DOCA TA usage.

```
#include <doca_model.p4>
#include <doca_headers.p4>
#include <doca_extens.p4>
#include <doca_parser.p4>
```

The DOCA TA features a single control, which requires a headers struct and standard metadata. User metadata and packet out metadata, defined by the user, are optional.

```
control hello_packet(  
  inout nv_headers_t headers,  
  in nv_standard_metadata_t std_meta,  
  inout nv_empty_metadata_t user_meta,  
  inout nv_empty_metadata_t pkt_out_meta  
) {  
  action drop() {  
    nv_drop();  
  }  
  action forward(bit<32> port) {  
    nv_send_to_port(port);  
  }  
  table forward_table {  
    key = {  
      headers.ethernet.dst_addr : exact;  
    }  
    actions = {  
      drop;  
      forward;  
      NoAction;  
    }  
    default_action = forward(3);  
    const entries = {  
      (48w0x001111111111) : forward(1);  
      (48w0x002222222222) : forward(2);  
      (48w0x00dddddddddd) : drop();  
      (48w0x00aaaaaaaaaa) : NoAction();  
      (48w0x00bbbbbbbbbb) : NoAction();  
    }  
  }  
  apply {  
    forward_table.apply();  
  }  
}
```

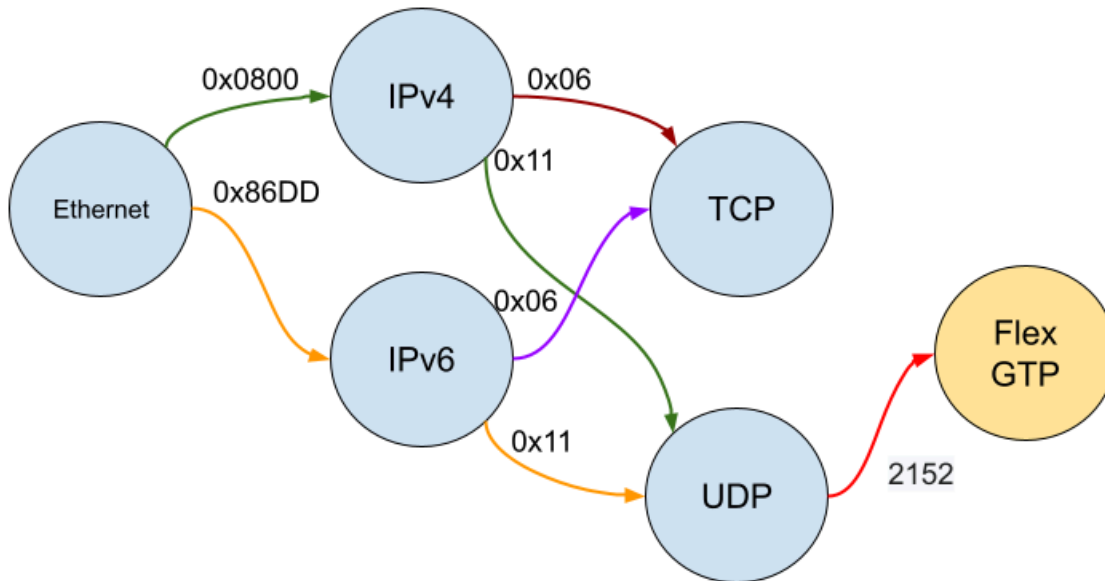
Finally, the main package must be instantiated:

```
NvDocaPipeline(  
  nv_fixed_parser(),  
  hello_packet()  
) main;
```

See the full DPL example [hello\\_packet.p4](#)

# GTP Parsing Example

This example demonstrates how to add a simple flex parser node to the existing hardware defined parse graph.



## Sample Code

The example starts with some basic definitions.

```
#include <doca_model.p4>
#include <doca_headers.p4>
#include <doca_externs.p4>
#include <doca_parser.p4>

const bit<32> WIRE_PORT = 32w00;
const bit<32> GTP_VPORT = 32w01;
const bit<32> DEFAULT_VPORT = 32w04;

struct metadata_t {
}
#define GTP_U_PORT 2152
```



Then we define the GTP-U version 1 header.

```
header Gtp_v1_h {
    bit<3>      version;          /** For GTPv1, this has a value of 1. */
    bit        protocol_type;    /** GTP (value 1) from GTP' (value 0) */
    bit        reserved;
    bit        extension_header_flag; /** extension header optional field. */
    bit        seq_number_flag;   /** Sequence Number optional field */
    bit        n_pdu_number_flag; /** N-PDU number optional field */
    bit<8>     message_type;      /** types of messages are defined in
3GPP TS 29.060 section 7.1 */
    bit<16>    message_length;    /** length of the payload in bytes */
    bit<32>    teid;              /** Tunnel endpoint identifier */
    bit<16>    sequence_number;   /** optional */
    bit<8>     n_pdu_number;      /** optional */
    bit<8>     next_extension_hdr_type; /** optional if any of the E, S, or PN
bits are on. The field must be interpreted only if the E bit is on */
}
```

Then we add NV\_FIXED\_HEADERS to the headers struct, along with the new GTP header.

```
struct headers_t {
    NV_FIXED_HEADERS
    Gtp_v1_h    gtpv1;
}
```

Using the `nv_transition_from` annotation, the GTP parser state is connected as a select transition from the UDP state.

```
parser packet_parser(packet_in packet, out headers_t headers) {
    NV_FIXED_PARSER(packet, headers)
```

```

@nv_transition_from("nv_parse_udp", GTP_U_PORT)
state parse_gtp
{
    packet.extract(headers.gtpv1);
    transition accept;
}
}

```

The control example uses a single flow table that matches on input port and GTP tunnel endpoint ID. The policy is then to forward the GTP packet to a port or drop the packet.

```

/**
 * This control admits GTP packets only if the tunnel ID matches
 *
 */
control gtp_tunnel(
    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES) gtp_counter;

    action send_to_port(nv_logical_port_t port) {
        gtp_counter.count();
        nv_send_to_port(port);
    }

    action drop() {
        gtp_counter.count();
        nv_drop();
    }

    table gtp_table {
        key = {

```

```

        std_meta.ingress_port: exact;
        headers.gtpv1.teid: exact;
    }
    actions = {
        send_to_port;
        drop;
    }
    default_action = drop;
    direct_counter = gtp_counter;

    const entries = {
        (WIRE_PORT, 0x00000001) : send_to_port(GTP_VPORT);
        (GTP_VPORT, 0x00000001) : send_to_port(WIRE_PORT);
    }
}

apply {
    if (headers.gtpv1.isValid()) {
        if (gtp_table.apply().miss) {
            nv_send_to_port(DEFAULT_VPORT);
        }
    }
    drop();
}
}

NvDocaPipeline(
    packet_parser(),
    gtp_tunnel()
) main;

```

See the full DPL example [gtp\\_parsing.p4](#)

## GTP Tunnel Encapsulation Example

The following example demonstrates GTPv1 encapsulation which must be declared as a custom header:

```
header gtp_v1_h {
    bit<3>      version;          /** For GTPv1, this has a value of 1. */
    bit        protocol_type;    /** GTP (value 1) from GTP' (value 0) */
    bit        reserved;
    bit        extension_header_flag; /** extension header optional field. */
    bit        seq_number_flag;   /** Sequence Number optional field */
    bit        n_pdu_number_flag; /** N-PDU number optional field */
    bit<8>     message_type;     /** types of messages are defined in
3GPP TS 29.060 section 7.1 */
    bit<16>    message_length;    /** length of the payload in bytes */
    bit<32>    teid;             /** Tunnel endpoint identifier */
    bit<16>    sequence_number;   /** optional */
    bit<8>     n_pdu_number;     /** optional */
    bit<8>     next_extension_hdr_type; /** optional if any of the E, S, or PN
bits are on. The field must be interpreted only if the E bit is on */
}
```

The custom GTPv1 header can optionally be added to the parse graph (see [Custom Parser States](#)) is optional, and is only required if any of the custom header fields need to be read from or modified after encapsulation.

```
struct headers_t {
    NV_FIXED_HEADERS
    gtp_v1_h    gtpv1;
}

parser packet_parser(packet_in packet, out headers_t headers) {
    NV_FIXED_PARSER(packet, headers)

    @nv_transition_from("nv_parse_udp", GTP_U_PORT)
    state parse_gtp
```

```

    {
        packet.extract(headers.gtpv1);
        transition nv_parse_inner_ipv4;
    }
}

```

The extern object `NvTunnelTemplate<HEADER_TYPE>` requires a struct type, for which the entire underlay header is defined as a struct:

```

struct tunnel_headers_t {
    nv_ethernet_h ethernet;
    nv_ipv4_h     ipv4;
    nv_udp_h      udp;
    gtp_v1_h      gtpv1;
}

```

Finally, when declaring a `NvTunnelTemplate` object, the annotation `@nv_tunnel_fields` must be present. This annotation contains a key-value entry for each header in the specified struct type. Under each header, each header field must be specified in the order that the field appears in the header. For the GTPv1 example, a `NvTunnelTemplate` would be declared in the main control:

```

@nv_tunnel_fields(
    ethernet = {
        dst_addr = "variable",
        src_addr = "variable",
        ether_type = 0x0800
    },
    ipv4 = {
        version = 0x4,
        ihl = 0x5,
        diffserv = 0,
        ecn = "ignore", // Cannot set
    }
)

```

```

    total_len = "ignore",          // reparse will calculate, cannot set
    identification = "ignore",    // Cannot set
    flags = 0,
    frag_offset = 0,
    ttl = 64,
    protocol = 17,
    hdr_checksum = "ignore",     // reparse will calculate, cannot set
    src_addr = "variable",
    dst_addr = "variable"
},
udp = {
    src_port = "ignore",         // reparse will set, cannot be set manually
    dst_port = "variable",
    length = "ignore",          // reparse will calculate, cannot set
    checksum = "ignore"         // Cannot set
},
gtpv1 = {
    version = 1,
    protocol_type = 1,
    reserved = 0,
    extension_header_flag = 1,
    seq_number_flag = 1,
    n_pdu_number_flag = 1,
    message_type = 0xFF,        // T-PDU
    message_length = "variable",
    teid = "variable",
    sequence_number = 0,
    n_pdu_number = 0,
    next_extension_hdr_type = 0
}
)
NvTunnelTemplate<tunnel_headers_t>() gtpv1Tunnel;

```

The type `tunnel_headers_t` determines the required structure of the annotation. Fields marked "variable" will be determined according to the arguments passed to

`nv_set_l2tunnel_underlay` or `nv_set_l3tunnel_underlay`. Some fields will always be overwritten by reparse, such as IPv4 length, IPv4 checksum, UDP length, etc which must be calculated and cannot be set by the user. Such fields can be marked "ignore". Marking other fields "ignore", would set them to zero.

Finally, the instantiated `NvTunnelTemplate` object can be used in one of the custom tunnel externs. For example:

```
action forward(nv_ipv4_addr_t src_addr) {
    nv_set_l3tunnel_underlay(headers, gtpv1Tunnel, {
        48w0x001A2B3C4D5E, // ethernet.dst_addr
        48w0x00F1E2D3C4B5, // ethernet.src_addr
        src_addr,           // ipv4.src_addr
        32w0xC07B7B01,     // ipv4.dst_addr
        16w0x868,          // udp.dst_port for GTP-U
        16w0x4,            // gtpv1.message_length
        32w0x1234567       // gtpv1.teid
    });
    nv_send_to_port(1);
}
```

The third argument to `nv_set_l3tunnel_underlay` is required to be a list expression where each element corresponds to a header field marked "variable" in the `NvTunnelTemplate` annotation. The order of the values in the list is the order in which the header fields appear in the packet.

[Full example of gtp\\_tunnel.p4:](#)

```
/*
 * SPDX-FileCopyrightText: Copyright (c) 2025 NVIDIA CORPORATION & AFFILIATES. All rights reserved.
 * SPDX-License-Identifier: LicenseRef-NvidiaProprietary
 *
 * NVIDIA CORPORATION, its affiliates and licensors retain all intellectual
 * property and proprietary rights in and to this material, related
 * documentation and any modifications thereto. Any use, reproduction,
 * disclosure or distribution of this material and related documentation
 * without an express license agreement from NVIDIA CORPORATION or
```

```
* its affiliates is strictly prohibited.
```

```
*/
```

```
#include <doca_model.p4>
```

```
#include <doca_headers.p4>
```

```
#include <doca_extens.p4>
```

```
#include <doca_parser.p4>
```

```
#define GTP_U_PORT 2152
```

```
header gtp_v1_h {
```

```
    bit<3>          version;                /** For GTPv1, this has a value of 1. */
```

```
    bit            protocol_type;          /** GTP (value 1) from GTP' (value 0) */
```

```
    bit            reserved;
```

```
    bit            extension_header_flag;  /** extension header optional field. */
```

```
    bit            seq_number_flag;       /** Sequence Number optional field */
```

```
    bit            n_pdu_number_flag;     /** N-PDU number optional field */
```

```
    bit<8>         message_type;          /** types of messages are defined in
```

```
3GPP TS 29.060 section 7.1 */
```

```
    bit<16>        message_length;        /** length of the payload in bytes */
```

```
    bit<32>        teid;                  /** Tunnel endpoint identifier */
```

```
    bit<16>        sequence_number;       /** optional */
```

```
    bit<8>         n_pdu_number;          /** optional */
```

```
    bit<8>         next_extension_hdr_type; /** optional if any of the E, S, or PN
```

```
bits are on. The field must be interpreted only if the E bit is on */
```

```
}
```

```
struct headers_t {
```

```
    NV_FIXED_HEADERS
```

```
    gtp_v1_h    gtpv1;
```

```
}
```

```
parser packet_parser(packet_in packet, out headers_t headers) {  
    NV_FIXED_PARSER(packet, headers)
```

```
    @nv_transition_from("nv_parse_udp", GTP_U_PORT)
```

```
    state parse_gtp
```



```

    {
        packet.extract(headers.gtpv1);
        transition nv_parse_inner_ipv4;
    }
}

struct tunnel_headers_t {
    nv_ethernet_h ethernet;
    nv_ipv4_h      ipv4;
    nv_udp_h       udp;
    gtp_v1_h       gtpv1;
}

control c(
    inout headers_t headers,
    in  nv_standard_metadata_t std_meta,
    inout nv_empty_metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    @nv_tunnel_fields(
        ethernet = {
            dst_addr = "variable",
            src_addr = "variable",
            ether_type = 0x0800
        },
        ipv4 = {
            version = 0x4,
            ihl = 0x5,
            diffserv = 0,
            ecn = "ignore", // cannot set
            total_len = "ignore", // HW will calculate, cannot set
            identification = "ignore", // cannot set
            flags = 0,
            frag_offset = 0,
            ttl = 64,
            protocol = 17,

```

```

        hdr_checksum = "ignore", // HW will calculate, cannot set
        src_addr = "variable",
        dst_addr = "variable"
    },
    udp = {
        src_port = "ignore", // HW will calculate entropy, cannot set
        dst_port = GTP_U_PORT,
        length = "ignore", // HW will calculate, cannot set
        checksum = "ignore" // cannot set
    },
    gtpv1 = {
        version = 1,
        protocol_type = 1,
        reserved = 0,
        extension_header_flag = 1,
        seq_number_flag = 1,
        n_pdu_number_flag = 1,
        message_type = 0xFF,
        message_length = "variable",
        teid = "variable",
        sequence_number = 123,
        n_pdu_number = 255,
        next_extension_hdr_type = 0
    }
)
NvTunnelTemplate<tunnel_headers_t>() gtpv1Tunnel;

action drop() {
    nv_drop();
}

action forward(nv_logical_port_t port) {
    nv_send_to_port(port);
}

action tunnel(nv_ipv4_addr_t dst_addr, bit<32> teid) {

```

```

    nv_set_l3tunnel_underlay(headers, gtpv1Tunnel, {
        48w0x1, 48w0x2, 32w0x3, dst_addr,
        16w4, /* 4 bytes of optional GTP-U payload */
        teid
    });
}

table ip_as_key {
    key = {
        headers.ipv4.src_addr : exact;
    }
    actions = {
        NoAction;
        drop;
        forward;
        tunnel;
    }
    size = 128;
    default_action = NoAction;
    const entries = {
        (32w0x01010101) : tunnel(32w0x11111111, 1);
        (32w0x02020202) : tunnel(32w0x22222222, 2);
        (32w0x03030303) : tunnel(32w0x33333333, 3);
    }
}

table teid_as_key {
    key = {
        headers.gtpv1.teid : exact;
    }
    actions = {
        NoAction;
        drop;
        forward;
    }
    size = 128;
}

```

```

    default_action = forward(3);
    const entries = {
        (32w0x01) : forward(1);
        (32w0x02) : forward(2);
    }
}

apply {
    ip_as_key.apply();
    if (headers.gtpv1.isValid()) {
        teid_as_key.apply();
    }
    forward(3);
}
}

NvDocaPipeline(
    packet_parser(),
    c()
) main;

```

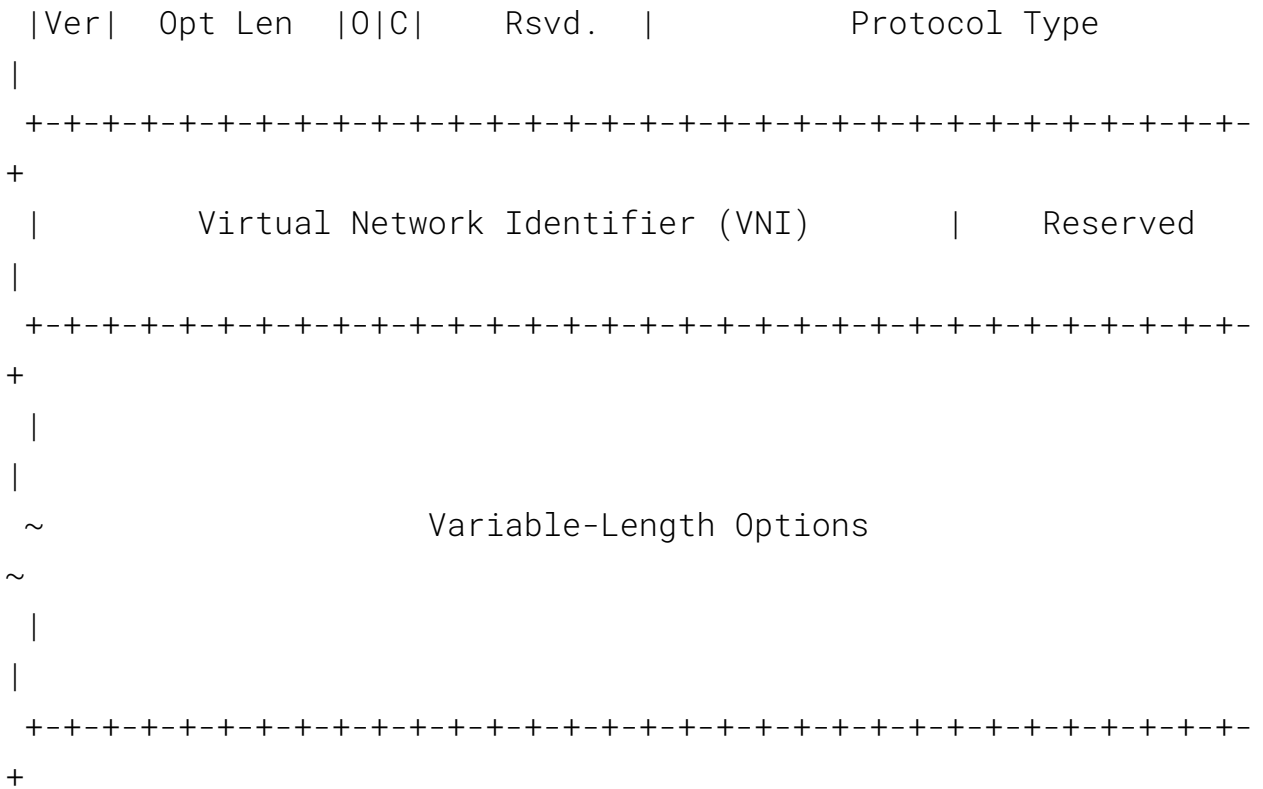
## Geneve TLV Parsing Example

This example demonstrates how to add a custom protocol header that supports optional Type-Length-Value (TLV) fields. In the P4-16 specification, there is no first class support in the parser primitives, and the user needs to manually.

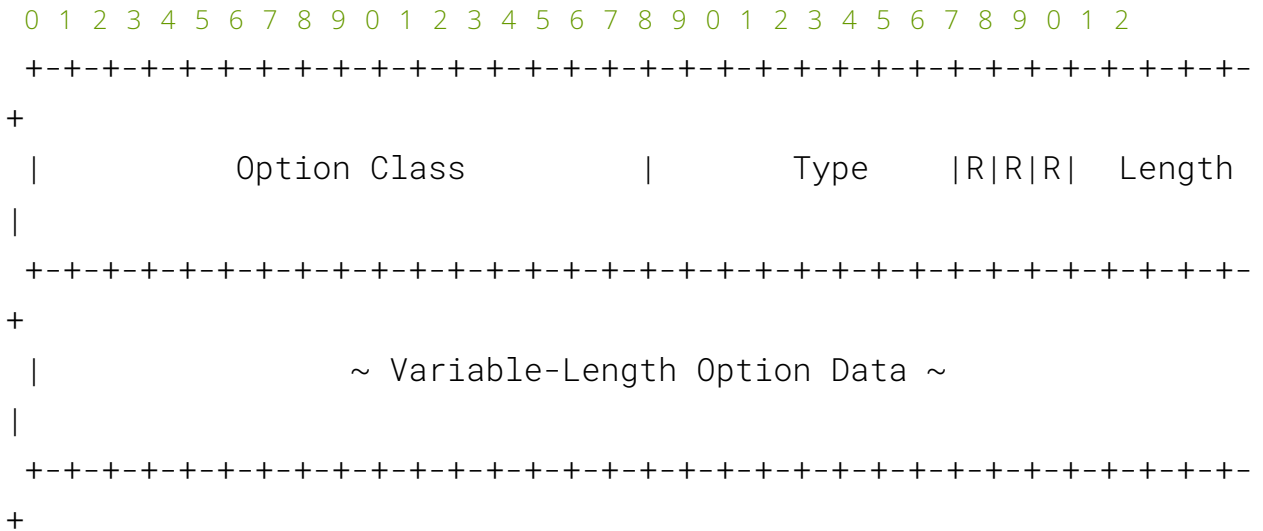
Geneve is short for Generic Network Virtualization Encapsulation, a common L2 tunnel header. The Geneve header format is:

```

 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
+
```



The Geneve option header is:



## Sample Code

First we define the Geneve header, per [RFC 8926](#).

```

header geneve_t {
    bit<2> ver;
    bit<6> opt_len;
    bit<1> o;
    bit<1> c;
    bit<6> reserved1;
    bit<16> protocol_type;
    bit<24> vni;
    bit<8> reserved2;
};

```

Next we define a struct that contains the base Geneve Option header, shared by all option types.

```

// Struct so it can be a field within other headers
struct geneve_option_t {
    bit<24> option_class_type;
    bit<3> reserved;
    bit<5> length;
};

```

Lastly, depending on the option class and type, the variable length option data can be defined. In this example, we model the Geneve options used by the [P4 Inband Telemetry \(INT\) specification](#).

```

// Intermediate nodes (INT Transit Hops) must process this type of INT Header.
header geneve_option_int_md_t {
    geneve_option_t base;
    bit<4> version;
    bit<1> discard;
    bit<1> exceeded_max_hops;
    bit<1> mtu_exceeded;
};

```

```

    bit<12> reserved;
    bit<5> hop_ml;
    bit<8> remaining_hop_count;
    bit<16> instruction_bitmap;
    bit<16> ds_id;
    bit<16> ds_instruction;
    bit<16> ds_flags;
};
// Destination headers can be used to enable Edge-to-Edge communication between
// the INT Source and INT Sink.
header geneve_option_int_destination_t {
    geneve_option_t base;
    // Destination headers must only be consumed by the INT Sink
};

// Intermediate nodes (INT Transit Hops) must process this type of INT Header
// and generate reports to the monitoring system as instructed.
header geneve_option_int_mx_t {
    geneve_option_t base;
    bit<4> version;
    bit<1> discard;
    bit<27> reserved1;
    bit<16> instruction_bitmap;
    bit<16> ds_id;
    bit<16> ds_instruction;
    bit<16> ds_flags;
};

```

All the above headers and structs must be added to the application's headers struct, along with the NV\_FIXED\_HEADERS.

```

struct app_headers {
    NV_FIXED_HEADERS

    geneve_t custom_geneve;
};

```

```

// geneve_option_t is not explicitly used by the P4 program, so a reference
// placed in the headers the type is not optimized out. The NvOptionParser
// instantiation references this type it by string name.
geneve_option_t geneve_opt;
geneve_option_int_md_t geneve_opt_int_md;
geneve_option_int_destination_t geneve_opt_int_destination;
geneve_option_int_mx_t geneve_opt_int_mx;
};

```

The last step for the parser is to define the parser state that will perform Geneve parsing. An extern object, NvOptionsParser, is used to further parse into the Geneve options. In this example, the TLV "type" is a combination of Option class and type:

- Class 0x103, type 1 (INT-MD)
- Class 0x103, type 2 (Destination-type)
- Class 0x103, type 3 (INT-MX)

```

parser geneve_parser(
    packet_in packet,
    out app_headers headers
) {
    NvOptionParser<bit<24>, _>(
        "opt_len",                // options_length_field
        2,                        // options_length_shift
        0,                        // options_length_add
        "geneve_option_t",        // option_layout_header_type
        "length",                // option_length_field
        0,                        // option_length_shift
        4,                        // option_length_add
        "option_class_type",      // option_type_field
        (list<tuple<bit<24>, _>>){ // options
            {24w0x010301, "headers.geneve_opt_int_md"},
            {24w0x010302, "headers.geneve_opt_int_destination"},
            {24w0x010303, "headers.geneve_opt_int_mx"}
        }
    );
}

```



```

    }
) geneveOptions;

NV_FIXED_PARSER(packet, headers)

@nv_transition_from("nv_parse_udp", GENEVE_PORT)
state parse_geneve {
    // Fixed geneve is only an example; "base" header must be flex.
    packet.extract(headers.custom_geneve);
    geneveOptions.parseOptions(packet, headers);
    transition select(headers.custom_geneve.protocol_type) {
        NV_TYPE_IPV4 : nv_parse_inner_ipv4;
        NV_TYPE_IPV6 : nv_parse_inner_ipv6;
        NV_TYPE_MAC : nv_parse_inner_ethernet;
        default : accept;
    }
}
}
}

```

The control below shows how a match action table can be configured to now match on a specific INT domain, and perform per ID forwarding.

```

control int_over_geneve(
    inout app_headers headers,
    in nv_standard_metadata_t std_meta,
    inout nv_empty_metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    action forward(bit<32> port) {
        nv_send_to_port(port);
    }

    action drop() {
        nv_drop();
    }
}

```

```

}

table geneve_option_int_md_table {
    key = {
        headers.geneve_opt_int_md.ds_id : exact;
    }
    actions = {
        forward;
        NoAction;
    }
    default_action = NoAction();
    const entries = {
        (16w0x100) : forward(1);
    }
}

table geneve_option_int_destination_table {
    key = {
        headers.custom_geneve.vni : exact;
    }
    actions = {
        forward;
        NoAction;
    }
    default_action = NoAction();
    const entries = {
        (24w0x200) : forward(2);
    }
}

table geneve_option_int_mx_table {
    key = {
        headers.geneve_opt_int_mx.ds_id : exact;
    }
    actions = {
        forward;
    }
}

```

```

        NoAction;
    }
    default_action = NoAction();
    const entries = {
        (16w0x300) : forward(3);
    }
}

apply {
    if (headers.custom_geneve.isValid()) {
        geneve_option_int_md_table.apply();
        geneve_option_int_destination_table.apply();
        geneve_option_int_mx_table.apply();
    }
    drop();
}
}

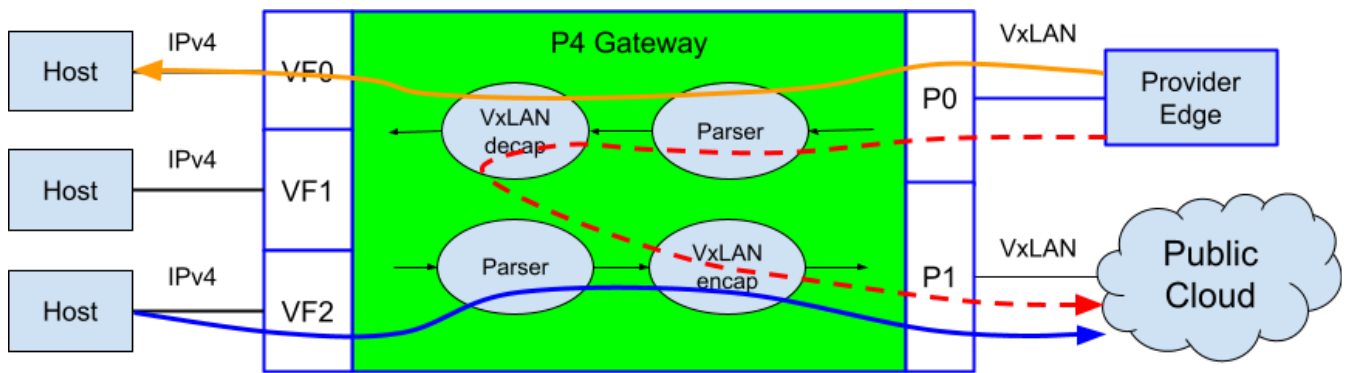
NvDocaPipeline(
    geneve_parser(),
    int_over_geneve()
) main;

```

See the full DPL example [geneve\\_tlv.p4](#)

## VXLAN Tunnel Gateway Example

This example illustrates how to write a basic VXLAN tunnel gateway using the DOCA target architecture. A tunnel gateway allows programmatic control over how VXLAN traffic can be "stitched" packets across tenant domains. In this example, end point traffic destined to local bare metal hosts can be decapsulated and forwarded to a VF, while gateway traffic can be decapsulated, re-encapsulated and sent back to the wire. For example, this program can be easily extended to be a gateway connecting legacy NVGRE networks to a VXLAN-GPE network.



## Sample Code

This example uses the native parser and 2 tables, of size 32K each. The wire port is configured with P4 port ID 0. A single bit in the user metadata structure is used to keep the decapsulation state.

```
#include <doca_model.p4>
#include <doca_headers.p4>
#include <doca_externs.p4>
#include <doca_parser.p4>

/*
 * Table sizes.
 */
const bit<32> DECAP_TABLE_SIZE = 32768;
const bit<32> ENCAP_TABLE_SIZE = 32768;

/* The directionality is based on network to host
 * The user will configure the P4 port IDs in the OVS configuration
 */
const bit<32> WIRE_PORT = 32w0;

struct metadata_t {
    bit<1> was_decapped;
}

struct headers_t {
    NV_FIXED_HEADERS
```

```

}

parser packet_parser(packet_in packet, out headers_t headers) {
    NV_FIXED_PARSER(packet, headers)
}

```

The encapsulation control has a single table, matching on IPv4 destination address. If an entry matches, the packet is VXLAN encapsulated and forwarded to the specified port. If the packet does not hit any entry, then the packet is dropped. It is simple to add more complex policy rules such as a 5 tuple ACL.

```

/*
 * This control performs the overlay policy including L2 encap with VXLAN
 */
control overlay_encap(
    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES)
encap_counter;

    action deny() {
        nv_drop();
    }

    action vxlan_v4_encap(nv_mac_addr_t underlay_src_mac,
nv_mac_addr_t underlay_dst_mac,
                        nv_ipv4_addr_t underlay_sip, nv_ipv4_addr_t
underlay_dip, bit<24> vni, nv_logical_port_t port) {
        encap_counter.count();
        nv_set_vxlan_v4_underlay(headers, false, underlay_dst_mac,
underlay_src_mac, 0, underlay_sip, underlay_dip, vni);
        nv_send_to_port(port);
    }
}

```

```

}
table encap_v4_table {
    key = {
        headers.ipv4.dst_addr : exact;
    }
    actions = {
        vxlan_v4_encap;
        deny;
    }
    size = ENCAP_TABLE_SIZE;
    default_action = deny;
    direct_counter = encap_counter;
}

apply {
    if (headers.ipv4.isValid() && (user_meta.was_decapped ==
1)) {
        encap_v4_table.apply();
    }
}
}

```

The decapsulation control simply checks if the packet is VXLAN, and decapsulates it. From there the packet can be sent directly to a port, or hair-pinned back to the wire.

```

/*
 * This control is for packets from wire to host (RX)
 * and includes policy for L2 decap
 */
control decap_flow(
    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {

```

```

    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES)
decap_counter;

    action deny() {
        nv_drop();
    }

    action decap() {
        decap_counter.count();
        nv_l2_decap(headers);
        user_meta.was_decapped = 1;
    }

    action to_port(nv_logical_port_t port) {
        nv_send_to_port(port);
    }

    action decap_to_port(nv_logical_port_t port) {
        decap_counter.count();
        user_meta.was_decapped = 1;
        nv_l2_decap(headers);
        nv_send_to_port(port);
    }

    table decap_v4_table {
        key = {
            headers.vxlan.vni : exact;
        }
        actions = {
            decap;
            to_port;
            decap_to_port;
            deny;
            NoAction;
        }
        size = DECAP_TABLE_SIZE;
    }

```

```

        direct_counter = decap_counter;
        default_action = deny;
    }

    apply {
        if (headers.vxlan.isValid()) {
            decap_v4_table.apply();
        }
    }
}

```

The main control checks the ingress port to determine if the packet is Network to Host, or Host to Network. Depending on the direction, it applies the decap\_flow control, or performs an overlay encapsulation.

```

control gateway(
    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    overlay_encap() over;
    decap_flow() decap;

    /* user should add entries that correspond to the wire ports
    * A hit means this is an RX packet, miss means a TX packet
    */
    table direction_table {
        key = {
            std_meta.ingress_port : exact;
        }
        actions = {
            NoAction;
        }
        default_action = NoAction;
    }
}

```



```

    const entries = {
        (WIRE_PORT) : NoAction();
    }
}

apply {
    user_meta.was_decapped = 0;
    if (direction_table.apply().hit) {
        decap.apply(headers, std_meta, user_meta,
pkt_out_meta);
    }
    over.apply(headers, std_meta, user_meta, pkt_out_meta);
}
}

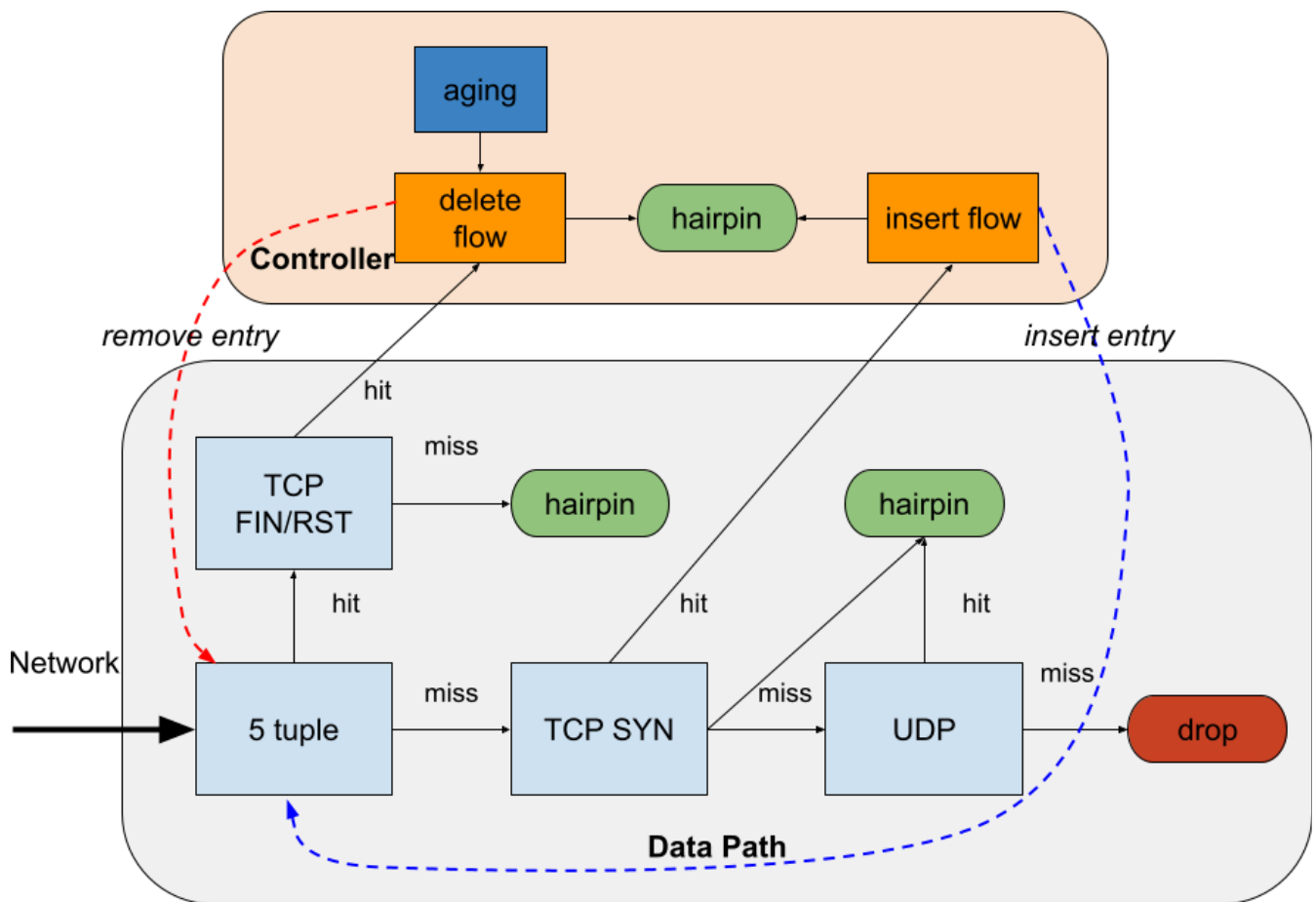
NvDocaPipeline(
    packet_parser(),
    gateway()
) main;

```

See the full DPL example [gateway.p4](#)

## Connection Tracking Example

Connection tracking is a stateful process that follows the a TCP/UDP session from establishment to termination. It consists of a 5 tuple match on <Source IP address, Destination IP address, IP protocol, L4 Source port, L4 Destination port>. The DPL datapath relies on a controller to insert entries for new 5 tuple flows that have not been seen yet.



## Sample Code

Connection metadata is a struct that sent with a packet that is sent to the P4Runtime controller. It requires a special annotation with the label "packet\_in". The user can define up to 32 bits of data in this data structure.

```
#include <doca_model.p4>
#include <doca_headers.p4>
#include <doca_externs.p4>
#include <doca_parser.p4>

@nv_controller_metadata("packet_in")
struct connection_meta_t {
    bit<2> type;
    bit<2> _reserved;
    bit<28> zone;
```

```
}
```

In this sample, the packet processing pipeline examines the state of the connection based on the TCP flags and set in the connection metadata an enum type describing the state of the flow.

```
/* Meta connection type */  
enum bit<2> ct_type {  
    PASS = 0,  
    NEW = 1,  
    RST = 2,  
    FIN = 3  
}
```

For readability, constants are used to assign P4 port IDs to more a more meaningful symbol name. The mapping of PFs and VFs to P4 port IDs is done separately in the [Service Configuration](#).

```
/*  
    VF port to SW entity that manages the flow insertion,  
    deletion and expiry times  
*/  
const bit<32> WIRE_PORT = 32w0;  
const bit<32> HAIRPIN_PORT = 32w1;
```

```
control conn_track(  
    inout nv_headers_t headers,  
    in nv_standard_metadata_t std_meta,  
    inout nv_empty_metadata_t user_meta,  
    inout nv_empty_metadata_t pkt_out_meta  
) {  
    connection_meta_t ctrl_meta;
```

```

/* 5-tuple matching for L4 TCP/UDP flows */
NvDirectCounter(NvCounterType.PACKETS_AND_BYTES)
table_t5_counter;

action set_connection_id_ct_table_t5(bit<28> zone) {
    table_t5_counter.count();
    ctrl_meta.zone = zone;
    ctrl_meta.type = ct_type.PASS;
}

action no_action_ct_table_t5() {
    table_t5_counter.count();
}

table ct_table_t5 {
    key = {
        headers.ipv4.src_addr : exact;
        headers.ipv4.dst_addr : exact;
        headers.ipv4.protocol : exact;
        headers.tcp.src_port : exact;
        headers.tcp.dst_port : exact;
    }
    actions = {
        set_connection_id_ct_table_t5;
        no_action_ct_table_t5;
    }
    direct_counter = table_t5_counter;
    // on hit table ct_table_known
    // on miss table ct_table_tcp_miss
    default_action = no_action_ct_table_t5;
    size = 1048576;
}

/*
* Known connections handling.

```

```

* Precondition - must be a TCP packet
* RST: TYPE_RST
* FIN: TYPE_FIN
* FINRST: TYPE_RST
*
* Match: tcp.flags
* Actions: meta.type, next PIPE
* MISS: next PIPE
*/
    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES)
table_known_counter;

    action set_connection_type_ct_table_known(bit<2> type) {
        table_known_counter.count();
        ctrl_meta.type = type;
    }

table ct_table_known {
    key = {
        headers.tcp.flags : exact;
    }
    actions = {
        set_connection_type_ct_table_known;
    }
    direct_counter = table_known_counter;
    default_action =
set_connection_type_ct_table_known(ct_type.PASS);
    const entries = {
        ( 0x1) :
set_connection_type_ct_table_known(ct_type.FIN);
        ( 0x4) :
set_connection_type_ct_table_known(ct_type.RST);
        (/*NV_TCP_PROTOCOL*/ 0x5) :
set_connection_type_ct_table_known(ct_type.RST); /* RST+FIN */
    }
}

```

```

    /*
    * Unknown TCP connections handling
    * Precondition - must be a TCP packet
    * SYN: TYPE_NEW
    *
    * Match: tcp.flags
    * Actions: PIPE, meta.type=NEW
    * MISS: NEXT PIPE
    */
    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES)
table_tcp_miss_counter;

    action set_connection_type_ct_table_tcp(bit<2> type) {
        table_tcp_miss_counter.count();
        ctrl_meta.type = type;
    }

    action no_action_ct_table_tcp() {
        table_tcp_miss_counter.count();
    }

    table ct_table_tcp {
        key = {
            headers.tcp.flags : exact;
        }
        actions = {
            set_connection_type_ct_table_tcp;
            no_action_ct_table_tcp;
        }
        direct_counter = table_tcp_miss_counter;
        default_action = no_action_ct_table_tcp;
        const entries = {
            (0x02) :
set_connection_type_ct_table_tcp(ct_type.NEW); /* SYN=1 ACK=0 */
        }
    }
}

```

```

    /*
    * Unknown UDP connections handling
    *
    * Match: UDP
    * Actions: PIPE, meta.type=NEW
    * MISS: DROP
    */
    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES)
table_udp_miss_counter;

    action set_connection_type_ct_table_udp(bit<2> type) {
        table_udp_miss_counter.count();
        ctrl_meta.type = type;
    }

    action drop_ct_table_udp() {
        table_udp_miss_counter.count();
        nv_drop();
    }

    table ct_table_udp {
        key = {
            std_meta.l4_type : exact;
        }
        actions = {
            set_connection_type_ct_table_udp;
            drop_ct_table_udp();
        }
        direct_counter = table_udp_miss_counter;
        default_action = drop_ct_table_udp();
        const entries = {
            (L4_TYPE_UDP) :
set_connection_type_ct_table_udp(ct_type.NEW); /* SYN=1 ACK=0 */
        }
    }

    /* user should add entries that correspond to the wire ports

```

\* A hit means this is an RX packet, miss means a TX packet

\*/

```
table direction_table {
    key = {
        std_meta.ingress_port : exact;
    }
    actions = {
        NoAction;
    }
    default_action = NoAction;
    const entries = {
        (WIRE_PORT) : NoAction();
    }
}

apply {
    ctrl_meta.type = ct_type.PASS;
    ctrl_meta._reserved = 0;
    ctrl_meta.zone = 0;

    if (direction_table.apply().hit) {
        if (std_meta.is_l4_ok == 1w1) {
            if (ct_table_t5.apply().hit) {
                if (ct_table_known.apply().hit) {
                    nv_send_to_controller(ctrl_meta);
                }
            } else if (headers.tcp.isValid()) {
                if (ct_table_tcp.apply().hit) {
                    nv_send_to_controller(ctrl_meta);
                }
            } else {
                if (ct_table_udp.apply().hit) {
                    nv_send_to_controller(ctrl_meta);
                }
            }
        }
    }
}
```



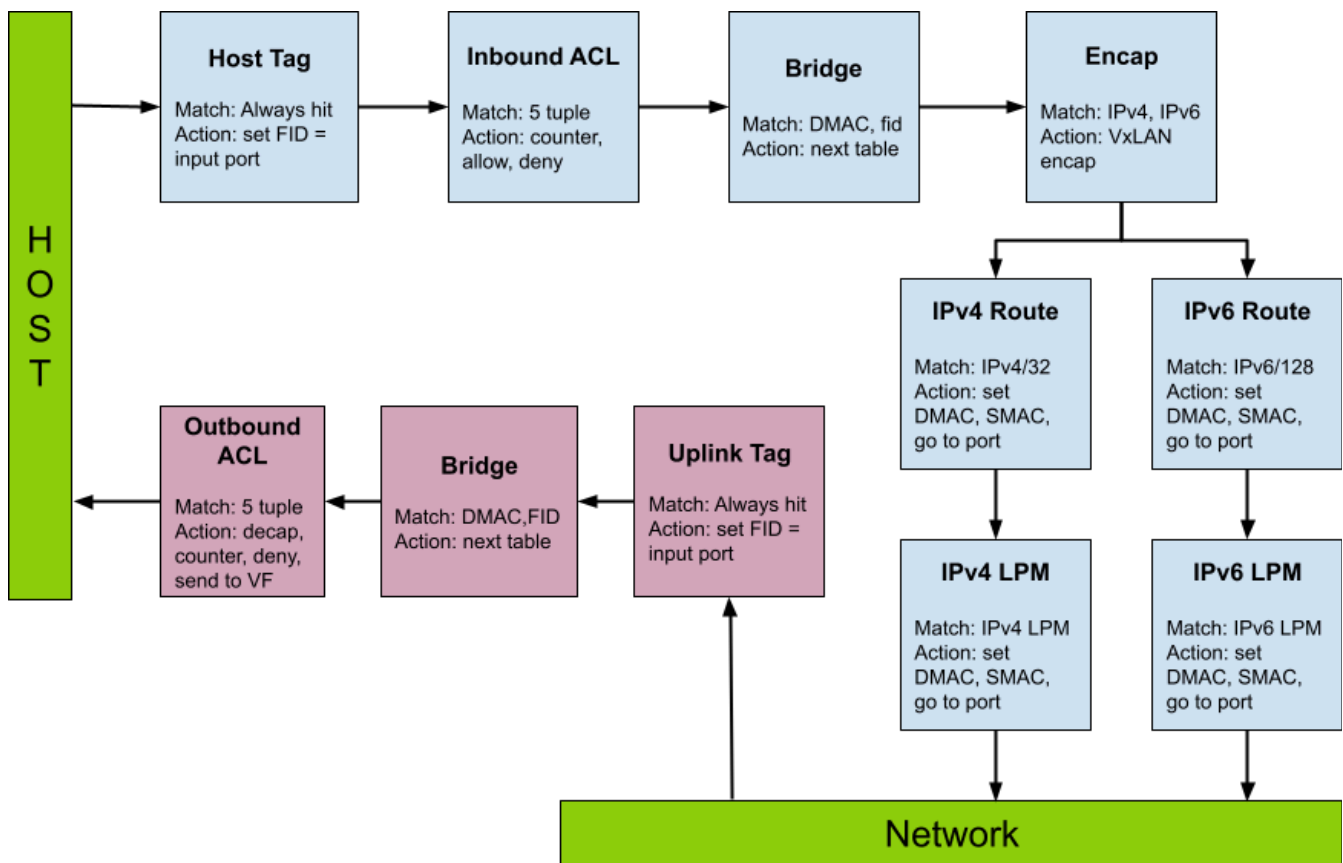
```
        }
        nv_send_to_port(HAIRPIN_PORT);
    }
}

// Instantiate the top-level DOCA Rx package
NvDocaPipeline(
    nv_fixed_parser(),
    conn_track()
) main;
```

See the full DPL example [conn\\_track.p4](#)

## Host-based Networking Example

HBN is service that provides controller-less VPC networking solution for bare metal as a service (BMaaS), featuring tenant isolation with EVPN VXLAN and accelerated routing on the host. This example demonstrates how to build a DPL datapath that supports host-based networking (HBN).



## Sample Code

This example uses indirect counters to monitor the number of packets that are allowed, denied and encapsulated/decapsulated. A user defined forwarding ID is used for determining how the bridge behaves.

```
enum bit<32> AdmissionCounter_t {
    DENY = 0,
    ALLOW = 1,
    VXLAN_V4_ENCAP = 2,
    VXLAN_V6_ENCAP = 3
}
```

```
enum bit<32> DecapFlowCounter_t {
    DENY = 0,
    ALLOW = 1,
}
```

```

/* The directionality is based on network to host
 * The user will configure the P4 port IDs in OVS
 */
const bit<32> WIRE_PORT = 32w0;

/* fid is 32 bits
 */
typedef bit<32> fid_t;

struct metadata_t {
    fid_t fid;
}

```

The default parser is used for this example:

```

struct headers_t {
    NV_FIXED_HEADERS
}

parser packet_parser(packet_in packet, out headers_t headers) {
    NV_FIXED_PARSER(packet, headers)
}

```

Finally, the main control body invokes three separate sub-controls:

- `overlay_encap` – defines the logic and policies for applying VXLAN encapsulations along with ACLs
- `underlay_route` – applies a custom IPv4/IPv6 routing tables with exact and LPM matching. The next hop is set, TTL decremented, and the packet is routed.

```

/**
 * This control performs the overlay policy including L2 encap with VxLAN
 */
control overlay_encap(

```

```

    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES) fid_counter;
    NvCounter(4, NvCounterType.PACKETS_AND_BYTES)
admission_counter;

    action set_fid(fid_t fid) {
        user_meta.fid = fid;
    }
    action deny() {
        admission_counter.count(AdmissionCounter_t.DENY);
        nv_drop();
    }
    action allow() {
        admission_counter.count(AdmissionCounter_t.ALLOW);
    }
    action vxlan_v4_encap(nv_mac_addr_t underlay_src_mac,
nv_mac_addr_t underlay_dst_mac,
        nv_ipv4_addr_t underlay_sip, nv_ipv4_addr_t
underlay_dip, bit<24> vni) {
        nv_set_vxlan_v4_underlay(headers, false, underlay_dst_mac,
underlay_src_mac, 0, underlay_sip, underlay_dip, vni);

admission_counter.count(AdmissionCounter_t.VXLAN_V4_ENCAP);
    }
    action vxlan_v6_encap(nv_mac_addr_t underlay_src_mac,
nv_mac_addr_t underlay_dst_mac,
        nv_ipv6_addr_t underlay_sip, nv_ipv6_addr_t
underlay_dip, bit<24> vni) {
        nv_set_vxlan_v6_underlay(headers, false, underlay_dst_mac,
underlay_src_mac, 0, underlay_sip, underlay_dip, vni);

admission_counter.count(AdmissionCounter_t.VXLAN_V6_ENCAP);

```

```

}
// A FID could be as simple as ingress port to FID mapping
// e.g. ports 1 and 2 are in FID 1, ports 3 and 4 are in FID 2
table fid_table {
    key = {
        std_meta.ingress_port : exact;
    }
    actions = {
        set_fid;
        NoAction;
    }
    size = FID_TABLE_SIZE;
    default_action = NoAction;
    direct_counter = fid_counter;
}

table admit_v4_table {
    key = {
        headers.ipv4.src_addr : exact;
        headers.ipv4.dst_addr : exact;
        headers.ipv4.protocol : exact;
        headers.tcp.src_port : exact;
        headers.tcp.dst_port : exact;
    }
    actions = {
        allow;
        deny;
        NoAction;
    }
    size = ADMIT_TABLE_SIZE;
    default_action = NoAction;
}

table admit_v6_table {
    key = {
        headers.ipv6.src_addr : exact;

```

```

        headers.ipv6.dst_addr : exact;
        headers.ipv6.next_header : exact;
        headers.tcp.src_port : exact;
        headers.tcp.dst_port : exact;
    }
    actions = {
        allow;
        deny;
        NoAction;
    }
    size = ADMIT_TABLE_SIZE;
    default_action = NoAction;
}

table bridge {
    key = {
        headers.ethernet.dst_addr : exact;
        user_meta.fid : exact;
    }
    actions = {
        NoAction;
    }
    size = BRIDGE_TABLE_SIZE;
    default_action = NoAction;
}

table encap_v4_table {
    key = {
        headers.ipv4.dst_addr : lpm;
    }
    actions = {
        vxlan_v4_encap;
    }
    size = ENCAP_TABLE_SIZE;
    default_action = vxlan_v4_encap(DEFAULT_DST_MAC,
DEFAULT_SRC_MAC,

```

```

                                                                    DEFAULT_SRC_IPV4,
DEFAULT_DST_IPV4, DEFAULT_VNI);
    }

    table encap_v6_table {
        key = {
            headers.ipv6.dst_addr : lpm;
        }

        actions = {
            vxlan_v6_encap;
        }
        size = ENCAP_TABLE_SIZE;
        default_action = vxlan_v6_encap(DEFAULT_DST_MAC,
DEFAULT_SRC_MAC,
                                                                    DEFAULT_SRC_IPV6,
DEFAULT_DST_IPV6, DEFAULT_VNI);

    }
    apply {
        user_meta.fid = 0;
        fid_table.apply();
        if (headers.ipv4.isValid()) {
            admit_v4_table.apply();
            if (bridge.apply().hit) {
                encap_v4_table.apply();
            }
        }
        else if (headers.ipv6.isValid()) {
            admit_v6_table.apply();
            if (bridge.apply().hit) {
                encap_v6_table.apply();
            }
        }
    }
}
}

```

```

/**
 * This control performs the underlay policy with routing
 */
control underlay_route(
    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    action set_port_and_route(nv_logical_port_t port,
nv_mac_addr_t src_mac, nv_mac_addr_t dst_mac) {
        headers.ethernet.src_addr = src_mac;
        headers.ethernet.dst_addr = dst_mac;
        nv_dec_ip_ttl(headers, 1);
        nv_send_to_port(port);
    }

    table exact_v4_route {
        key = {
            headers.ipv4.dst_addr : exact;
        }
        actions = {
            set_port_and_route;
            NoAction;
        }
        size = ROUTE_TABLE_SIZE;
        default_action = NoAction;
    }

    table lpm_v4_route {
        key = {
            headers.ipv4.dst_addr : lpm;
        }
        actions = {
            set_port_and_route;
            NoAction;
        }
    }
}

```



```

    }
    size = ROUTE_TABLE_SIZE;
    default_action = NoAction;
}
table exact_v6_route {
    key = {
        headers.ipv6.dst_addr : exact;
    }
    actions = {
        set_port_and_route;
        NoAction;
    }
    size = ROUTE_TABLE_SIZE;
    default_action = NoAction;
}
table lpm_v6_route {
    key = {
        headers.ipv6.dst_addr : lpm;
    }
    actions = {
        set_port_and_route;
        NoAction;
    }
    size = ROUTE_TABLE_SIZE;
    default_action = NoAction;
}

apply {
    if (headers.ipv4.isValid()) {
        if (exact_v4_route.apply().miss) {
            lpm_v4_route.apply();
        }
    }
    else if (headers.ipv6.isValid()) {
        if (exact_v6_route.apply().miss) {
            lpm_v6_route.apply();
        }
    }
}

```

```

    }
  }
}

```

- `decap_flow` – defines the policies for forwarding ID, L2 bridging, decapsulation of VXLAN and VXLAN-GPE tunnels, and forwarding to the VF

```

/**
 * This control is for packets from wire to host (RX)
 * and includes policy for L2 decap
 */
control decap_flow(
    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    NvDirectCounter(NvCounterType.PACKETS_AND_BYTES) fid_counter;
    NvCounter(2, NvCounterType.PACKETS_AND_BYTES) decap_counter;

    action deny() {
        decap_counter.count(DecapFlowCounter_t.DENY);
        nv_drop();
    }
    action set_fid(fid_t fid) {
        user_meta.fid = fid;
    }
    action decap_l2_and_send(nv_logical_port_t port) {
        decap_counter.count(DecapFlowCounter_t.ALLOW);
        nv_l2_decap(headers);
        nv_send_to_port(port);
    }
    action decap_ipv4_and_send(nv_logical_port_t port) {
        decap_counter.count(DecapFlowCounter_t.ALLOW);
    }
}

```

```

        nv_l3_decap(headers, false, 0xffffffff,
0x112233445566, NV_TYPE_IPV4, 0);
        nv_send_to_port(port);
    }
    action decap_ipv6_and_send(nv_logical_port_t port) {
        decap_counter.count(DecapFlowCounter_t.ALLOW);
        nv_l3_decap(headers, false, 0xffffffff,
0x112233445566, NV_TYPE_IPV6, 0);
        nv_send_to_port(port);
    }

table fid_table {
    key = {
        std_meta.ingress_port : exact;
    }
    actions = {
        set_fid;
        NoAction;
    }
    size = FID_TABLE_SIZE;
    default_action = NoAction;
    direct_counter = fid_counter;
}

table bridge {
    key = {
        headers.ethernet.dst_addr : exact;
        user_meta.fid : exact;
    }
    actions = {
        deny;
        NoAction;
    }
    size = BRIDGE_TABLE_SIZE;
    default_action = deny;
}

```

```

table decap_v4_table {
    key = {
        headers.ipv4.src_addr : exact;
        headers.ipv4.dst_addr : exact;
        headers.ipv4.protocol : exact;
        headers.udp.src_port : exact;
        headers.udp.dst_port : exact;
    }
    actions = {
        decap_l2_and_send;
        decap_ipv4_and_send;
        deny;
    }
    size = DECAP_TABLE_SIZE;
    default_action = deny;
}

table decap_v6_table {
    key = {
        headers.ipv6.src_addr : exact;
        headers.ipv6.dst_addr : exact;
        headers.ipv6.next_header : exact;
        headers.udp.src_port : exact;
        headers.udp.dst_port : exact;
    }
    actions = {
        decap_l2_and_send;
        decap_ipv6_and_send;
        deny;
    }
    size = DECAP_TABLE_SIZE;
    default_action = deny;
}

apply {

```

```

user_meta.fid = (fid_t) 0;
fid_table.apply();
if (bridge.apply().hit) {
    if (headers.ipv4.isValid()) {
        decap_v4_table.apply();
    }
    else if (headers.ipv6.isValid()) {
        decap_v6_table.apply();
    }
}
}
}

```

The main control checks which direction the packet came from.

- if the packet is from network to host, it invokes the decap\_flow control
- if the packet is from host to network, it invokes the overlay\_encap and routing controls

```

control host_based_networking(
    inout headers_t headers,
    in nv_standard_metadata_t std_meta,
    inout metadata_t user_meta,
    inout nv_empty_metadata_t pkt_out_meta
) {
    overlay_encap() overlay;
    underlay_route() route;
    decap_flow() decap;

    /* user should add entries that correspond to the wire ports
    * A hit means this is an RX packet, miss means a TX packet
    */
    table direction_table {
        key = {

```

```

        std_meta.ingress_port : exact;
    }
    actions = {
        NoAction;
    }
    default_action = NoAction;
    const entries = {
        (WIRE_PORT) : NoAction();
    }
}

apply {
    if (direction_table.apply().miss) {
        overlay.apply(headers, std_meta, user_meta,
pkt_out_meta);
        route.apply(headers, std_meta, user_meta,
pkt_out_meta);
    }
    else {
        decap.apply(headers, std_meta, user_meta,
pkt_out_meta);
    }
}

}

NvDocaPipeline(
    packet_parser(),
    host_based_networking()
) main;

```

See the full DPL example [hbn.p4](#)

**Notice**  
This document is provided for information purposes only and shall not be regarded as a warranty of a certain functionality, condition, or quality of a product. NVIDIA Corporation (“NVIDIA”) makes no representations or warranties, expressed or implied, as to the accuracy or completeness of the information contained in this document and assumes no responsibility for any errors contained herein. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This document is not a commitment to develop, release, or deliver any Material (defined below), code,

or functionality.<br/><br/>NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and any other changes to this document, at any time without notice.<br/><br/>Customer should obtain the latest relevant information before placing orders and should verify that such information is current and complete.<br/><br/><br/>NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer (“Terms of Sale”). NVIDIA hereby expressly objects to applying any customer general terms and conditions with regards to the purchase of the NVIDIA product referenced in this document. No contractual obligations are formed either directly or indirectly by this document.<br/><br/>NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death, or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer’s own risk.<br/><br/>NVIDIA makes no representation or warranty that products based on this document will be suitable for any specified use. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer’s sole responsibility to evaluate and determine the applicability of any information contained in this document, ensure the product is suitable and fit for the application planned by customer, and perform the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer’s product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this document. NVIDIA accepts no liability related to any default, damage, costs, or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this document or (ii) customer product designs.<br/><br/>No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA.<br/><br/><br/><br/>Reproduction of information in this document is permissible only if approved in advance by NVIDIA in writing, reproduced without alteration and in full compliance with all applicable export laws and regulations, and accompanied by all associated conditions, limitations, and notices.<br/><br/><br/><br/>THIS DOCUMENT AND ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, “MATERIALS”) ARE BEING PROVIDED “AS IS.” NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL NVIDIA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF NVIDIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA’s aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the Terms of Sale for the product.<br/><br/><br/><br/><b>Trademarks</b><br/><br/><br/>NVIDIA and the NVIDIA logo are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.<br/>

© Copyright 2025, NVIDIA. PDF Generated on 04/24/2025