# Graph Specification

# Table of contents

Graph Specification is a format to describe high-performance AI applications in a modular and extensible way. It allows writing applications in a standard format and sharing components across multiple applications without code modification. Graph Specification is based on entity-composition pattern. Every object in graph is represented with entity (aka Node) and components. Developers implement custom components which can be added to entity to achieve the required functionality.

# Concepts

The graph contains nodes which follow an entity-component design pattern implementing the "composition over inheritance" paradigm. A node itself is just a light-weight object which owns components. Components define how a node interacts with the rest of the applications. For example, nodes be connected to pass data between each other. A special component, called compute component, is used to execute the code based on certain rules. Typically a compute component would receive data, execute some computation and publish data.

## Graph

A graph is a data-driven representation of an AI application. Implementing an application by using programming code to create and link objects results in a monolithic and hard to maintain program. Instead a graph object is used to structure an application. The graph can be created using specialized tools and it can be analyzed to identify potential problems or performance bottlenecks. The graph is loaded by the graph runtime to be executed.

The functional blocks of a graph are defined by the set of nodes which the graph owns. Nodes can be queried via the graph using certain query functions. For example, it is possible to search for a node by its name.

## SubGraph

A subgraph is a graph with additional node for interfaces. It points to the components which are accessible outside this graph. In order to use a subgraph in an existing graph or subgraph, the developer needs to create an entity where a component of the type `nvidia::gxf::Subgraph` is contained. Inside the Subgraph component a corresponding subgraph can be loaded from the `yaml` file indicated by *location* property and instantiated in the parent graph.

System makes the components from interface available to the parent graph when a sub-graph is loaded in the parent graph. It allows users to link sub-graphs in parent with defined interface.

A subgraph interface can be defined as follows:

> --- interfaces: - name: iname # the name of the interface for the access from the parent graph target: n_entity/n_component # the true component in the subgraph that is represented by the interface

## Node

Graph Specification uses an entity-component design principle for nodes. This means that a node is a light-weight object whose main purpose is to own components. A node is a composition of components. Every component is in exactly one node. In order to customize a node a developer does not derive from node as a base class, but instead composes objects out of components. Components can be used to provide a rich set of functionality to a node and thus to an application.

## Components

Components are the main functional blocks of an application. Graph runtime provides a couple of components which implement features like properties, code execution, rules and message passing. It also allows a developer to extend the runtime by injecting her own custom components with custom features to fit a specific use case.

The most common component is a codelet or compute component which is used for data processing and code execution. To implement a custom codelet you'll need to implement a certain set of functions like *start* and *stop*. A special system - the *scheduler* - will call these functions at the specified time. Typical examples of triggering code execution are: receiving a new message from another node, or performing work on a regular schedule based on a time trigger.

## Edges

Nodes can receive data from other nodes by connecting them with an edge. This essential feature allows a graph to represent a compute pipeline or a complicated AI application. An input to a node is called sink while an output is called source. There can be zero, one or multiple inputs and outputs. A source can be connected to multiple sinks and a sink can be connected to multiple sources.

**Extension**

An extension is a compiled shared library of a logical group of component type definitions and their implementations along with any other asset files that are required for execution of the components. Some examples of asset files are model files, shared libraries that the extension library links to and hence required to run, header and development files that enable development of additional components and extensions that use components from the extension.

An extension library is a runtime loadable module compiled with component information in a standard format that allows the graph runtime to load the extension and retrieve further information from it to:

- Allow the runtime to create components using the component types in the extension.

- Query information regarding the component types in the extension:

    - The component type name

    - The base type of the component

    - A string description of the component

    - Information of parameters of the component – parameter name, type, description etc.,

- Query information regarding the extension itself - Name of the extension, version, license, author and a string description of the extension.

The section :doc: *GraphComposer_Dev_Workflow* talks more about this with a focus on developing extensions and components.

# Graph File Format

Graph file stores list of entities. Each entity has a unique name and list of components. Each component has a name, a type and properties. Properties are stored as key-value pairs.

```
%YAML 1.2 --- name: source components: - name: signal type: sample::test::ping -
type: nvidia::gxf::CountSchedulingTerm parameters: count: 10 --- components: -
type: nvidia::gxf::GreedyScheduler parameters: realtime: false max_duration_ms:
1000000
```