

Laminar: A New Serverless Stream-based Framework with Semantic Code Search and Code Completion

Zaynab Zahra
University of St. Andrews
St Andrews, UK
zz46@st-andrews.ac.uk

Zihao Li
University of St. Andrews
St Andrews, UK
zl81@st-andrews.ac.uk

Rosa Filgueira
University of St. Andrews
St Andrews, UK
rf208@st-andrews.ac.uk

ABSTRACT

This paper introduces Laminar, a novel serverless framework based on `dispel4py`, a parallel stream-based dataflow library. Laminar efficiently manages streaming workflows and components through a dedicated registry, offering a seamless serverless experience. Leveraging large language models, Laminar enhances the framework with semantic code search, code summarization, and code completion. This contribution enhances serverless computing by simplifying the execution of streaming computations, managing data streams more efficiently, and offering a valuable tool for both researchers and practitioners.

KEYWORDS

serverless computing, streaming applications, semantic code search, transformers, `dispel4py`, code completion, code summarization

ACM Reference Format:

Zaynab Zahra, Zihao Li, and Rosa Filgueira. 2023. Laminar: A New Serverless Stream-based Framework with Semantic Code Search and Code Completion. In *Proceedings of 18th Workshop on Workflows in Support of Large-Scale Science (WORKS 2023)*. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

In the rapidly evolving landscape of cloud computing, serverless computing [15] has emerged as a transformative paradigm, offering scalability, cost-effectiveness, and simplicity in deploying applications. However, the surge in data-intensive applications and the increasing demand for real-time processing present new challenges [23] for existing serverless frameworks. Firstly, traditional serverless architectures struggle to efficiently handle the continuous flow of streaming data, resulting in bottlenecks and latency issues. Secondly, supporting stateful computations within a serverless environment becomes complex due to the need to maintain and manage the state across distributed and ephemeral instances.

To address these challenges we introduce Laminar¹, a novel open-source Serverless Stream-based Processing Framework with

¹<https://github.com/dispel4pyserverless>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WORKS 2023, November 12 2023, Denver, CO

© 2023 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Deep Learning Code Search. Unlike traditional serverless frameworks, Laminar provides a comprehensive solution that seamlessly handles data streams and supports stateful computations by leveraging the power of `dispel4py` Python library. `dispel4py` inherent support for parallelism enables concurrent data processing, while abstract workflow descriptions in Python empower users to construct intricate stream processing pipelines. Similar to how functions encapsulate specific tasks in traditional programming, `dispel4py` Processing Elements (PEs) represent modular computational units within the Laminar serverless environment. Furthermore, Laminar goes beyond existing serverless frameworks by providing novel deep learning search facilities to find relevant PEs. The main contributions of this work are:

- **Handling Data Streams and Stateful Computations:** Laminar bridges a gap in conventional serverless frameworks, providing native support for data streams and stateful computations through an intuitive workflow description system and a streamlined registry.
- **Endpoints:** Laminar implements comprehensive server-client architecture with endpoints to facilitate communication for registering, executing, and managing PEs and workflows.
- **Intuitive Client Functionality:** Laminar introduces an intuitive client interface for convenient PE and workflow registration, execution, and management. It automates library detection and passes the information to the server for remote PE execution.
- **Serverless Execution Engine:** Leveraging the principles of serverless computing, Laminar efficiently handles PE and workflow execution by automatically provisioning resources and installing the necessary libraries ensuring seamless serverless operation.
- **Registry:** Laminar provides a robust functionality for registering workflows, including PEs and their properties. The registry serves as a central repository, enhancing Laminar's efficiency and management of serverless components.
- **Deep Learning Code Search Facilities:** Laminar harnesses the power of large language models, enhancing its capabilities for advanced PE code search, code summarization, and code completion. This integration significantly enhances semantic code search functionality. We conducted evaluations of multiple models to determine the optimal ones for integration into Laminar.

The remainder of the paper is structured as follows. Section 2 presents background on technologies relevant for this work. Section 3 details the features of Laminar. While Section 4 specifically focuses on the different registry searches supported by Laminar.

The practical utility of Laminar is demonstrated through two showcased use cases in Section 5. Section 6 presents different evaluations performed with this framework. We contextualize our work by surveying related studies in Section 7, and finally, we conclude and outline future directions in Section 8.

2 BACKGROUND

In the following sections, we look into the background work that is inherently established in this paper.

2.1 dispel4py

dispel4py [4, 18] is a parallel stream-based dataflow framework designed for creating and executing data-intensive applications. It provides an abstract and user-friendly approach to workflow creation, including automatic parallelization, which allows users to design, execute, and optimize workflows without requiring in-depth knowledge of the underlying hardware or middleware. Key concepts in dispel4py include:

- **Processing Element (PE):** The fundamental unit in dispel4py, acting as a computational task within a workflow graph. PEs connect through inputs and outputs for seamless stream-based data flow. They can be stateful, retaining previous inputs, or stateless, focusing on current data. Various PE types are available for distinct roles: `ProducerPE` (one output port), `IterativePE` (one input port, one output port), `ConsumerPE` (one input port), `GenericPE` (custom-defined, any number of ports).
- **Instance:** Refers to a PE execution within the computing process. Multiple instances of a single PE can be assigned, allowing workflow scaling.
- **Connection:** Defines data flow between PEs, determining data consumption and production rates.
- **Mappings:** Maps workflows onto execution systems. These include a *Simple* mapping for running workflows sequentially, and parallel options such as *MPI* [6], *Redis* [2], and *Multiprocessing (Multi)*², eliminating the need for manual workflow modifications.
- **Abstract Workflow:** Represents logical connections between PEs, outlining computational sequences and data transformations. It is what the user describes.
- **Concrete Workflow:** During enactment (after the user specifies the mapping to use and the number of processes), dispel4py automatically builds the concrete (parallel³) workflow, which is a directed acyclic graph (DAG) based on the abstract workflow. The concrete workflow is the actual workflow executed by the compute infrastructure.
- **Grouping:** Specifies how PEs communicate during input connections. For example, `group-by` (see Listing 2), which behaves similarly to ‘MapReduce’, routing data units with the same value in the specified element to the same PE instance.

To create dispel4py workflows, users implement PEs and connect them in graphs. In Figure 1, an example workflow illustrates

three PEs distributed among five processes (e.g., one PE instance for PE1 and two for PE2 to PE3) using the *Multi* mapping. Users design the green abstract workflow, while the blue concrete workflow is automatically generated during enactment.

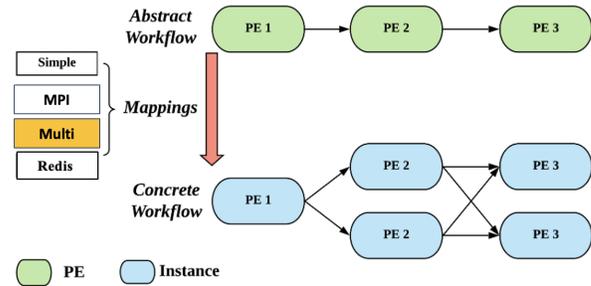


Figure 1: Example of a dispel4py workflow, in which the user has indicated to run it using the *Multi* mapping with five processes. Each PE instance runs in a different process.

Listing 1 provides the code for the first Processing Element (PE1) in the dispel4py workflow shown in Figure 1. Full workflow code is available in Listing 3. Note that the core functionality of a PE is encapsulated within the `_process` function.

```
class NumberProducer(ProducerPE):
    def __init__(self):
        ProducerPE.__init__(self)
    def _process(self):
        # Generate a random number
        result= random.randint(1, 1000)
        # Return the number as the output
        return result
```

Listing 1: NumberProducer *stateless* PE, referred to as PE1 in Figure 1, generates random numbers and streams them out.

2.2 Serverless Computing

Serverless computing [22], also known as Function-as-a-Service (FaaS), is a cloud computing paradigm that abstracts away server management and infrastructure concerns from developers. In a serverless architecture, developers focus solely on writing and deploying individual functions or code snippets to the cloud, and the cloud provider takes care of automatically provisioning, scaling, and managing the underlying infrastructure needed to execute those functions. This approach allows developers to focus on the core logic of their applications, without worrying about server maintenance, resource management, or scalability.

In a traditional cloud computing model, developers typically need to manage virtual machines (VMs) or containers to run their applications. This can be time-consuming and resource-intensive, especially when handling fluctuating workloads. Serverless computing, on the other hand, abstracts away the concept of servers

²<https://docs.python.org/3/library/multiprocessing.html>

³All mappings, except *Simple*, trigger the generation of an automatic parallel concrete workflow. This concrete workflow is capable of accommodating multiple parallel patterns, including the *farm* and *pipeline* patterns.

entirely, allowing developers to run their functions on demand without the need to explicitly manage the infrastructure. Functions are executed in stateless containers that are spun up and scaled automatically, based on the incoming workload.

One of the key advantages of serverless computing is its cost-effectiveness. Users are billed only for the actual execution time of their functions, rather than for idle server time. This pay-per-invocation pricing model can result in significant cost savings, especially for applications with variable workloads.

2.3 Language Models and Transformers

The landscape of computer capabilities in comprehending human languages, speaking, and reading has undergone a transformation due to the advent of natural language processing models. Among the state-of-the-art models, the transformer architecture has demonstrated remarkable advancements [28]. However, the scope of these models extends beyond human language and can be expanded to incorporate abstract syntax trees, enabling them to understand and compare code. In this work, we have applied three different language models:

- *UnixCoder*[8] is a specialized transformer-based model designed to convert Abstract Syntax Trees (ASTs) into sequential text representations [1]. This model enhances the semantic representation of code fragments through embeddings, achieved by employing multi-modal contrastive learning (MCL) for comprehensive code semantic capture using ASTs, and cross-modal generation (CMG) to align embeddings across different programming languages via code comments. In extensive experiments [8], *UnixCoder* outperformed state-of-the-art models like CodeBERT [3], GraphCodeBERT [9] and SYNCOBERT [25] in various code-related understanding tasks, including semantic code search and clone detection.
- *ReACC-py-retriever* is a model developed for *ReACC* retrieval-augmented code completion Framework [19]. The model utilizes “external” context for the code completion task by retrieving semantically and lexically similar codes from existing codebase. It has been pre-trained with a dual-encoder as a retriever for partial code search, which retrieves code fragments given a partial code. On the CodeXGLUE [20] benchmark, the model achieves a state-of-the-art performance in the code completion task.
- *CodeT5* [26] is a pre-trained encoder-decoder model that incorporates the token type information from code and allow for multi-task learning on downstream tasks. *CodeT5* is with three code intelligence capabilities: 1) *text-to-code generation* to generate code based on the natural language description ; 2) *Code autocompletion*, to complete the whole function of code given the target function name; and 3) *Code summarization* to generate the summary of a function in natural language description. Among those tasks, *CodeT5* achieves state-of-the-art performance for code summarization task.

2.4 Semantic Code Search Paradigms

Semantic code search [7] is a critical functionality in Laminar, enabling *text-to-code* similarity queries over the registry. This NLP

technique allows searching for existing code snippets through natural language, which can greatly improve programming efficiency. The objective is to identify matching codes in the search corpus that correspond to the query.

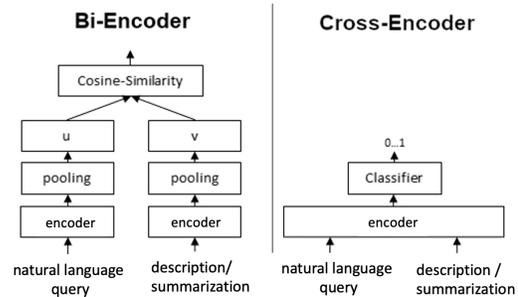


Figure 2: The concept diagram of *bi-encoder* (left) and *cross-encoder* (right) code search architecture.

We adopted the bi-encoder paradigm, as illustrated in Figure 2, where each input (natural language query or a code description/ summarization) is independently mapped into a dense vector space. Bi-encoders calculate embeddings for both inputs, enabling efficient storage of embeddings for subsequent queries. In contrast, cross-encoders perform full-attention over the input pairs of sentences, resulting in better accuracy but reduced efficiency. The bi-encoder architecture in Laminar is well-suited for applications like PEs similarity, where fast querying of separate embeddings is crucial. While bi-encoders are faster, cross-encoders achieve better accuracy but may not be practical in certain scenarios. In summary, the bi-encoder approach strikes a balance between efficiency and effectiveness, making it an ideal choice for searches in Laminar.

2.5 Code Completion and Code Summarization

In the context of this work, *Code completion* [24] involves predicting subsequent code tokens based on the given code context, contributing to enhanced programming productivity. Recent advancements have demonstrated the efficacy of statistical language modeling with transformers in improving code completion performance by leveraging vast source code datasets.

Conversely, code summarization [29] refers to the generation of concise and coherent natural language summaries or descriptions that encapsulate the core functionality and behavior of a specific segment of source code. This task aims to facilitate comprehension and documentation by providing human-readable explanations for intricate code snippets, functions, or methods.

In the Laminar framework, code completion serves the purpose of allowing users to input incomplete PE snippets for completion queries. The framework then retrieves similar PEs from our registry, which functions as our search corpus. Moreover, code summarization is leveraged for storing descriptions of PEs within the registry when users fail to provide them. This dual approach enhances the overall usability and efficiency of the framework’s capabilities.

clone-detection, as presented in [8]. Although the base model provided by the authors through Hugging Face ⁴ lacks fine-tuning tasks, we undertook our own fine-tuning process following instructions outlined in ⁵. This process involved utilizing the AdvTest dataset [20], which comprises 280,634 pairs of (documentation, function) sourced from CodeSearchNet [12]. Notably, the dataset normalizes Python function and variable names to enhance model understanding and generalization capabilities. Our fine-tuning efforts led to the development of two models: *unixcoder-code-search* and *unixcoder-clone-detection*. It is important to highlight that the fine-tuned model was originally developed by the authors of this work for another complementary study [17], which focused on comparing repository similarities. Each of these models underwent approximately 6 hours of fine-tuning on an NVIDIA A40 GPU server.

3 LAMINAR OVERVIEW

This section presents an overview of the fundamental components that constitute the Laminar architecture, each serving a distinct purpose within the framework. The core elements encompass the *Registry*, *Server*, *Execution Engine*, and *Client*, as detailed in Table 1.

Element	Purpose	Section
Registry	Stores user, PE, and workflow information	Section 3.1.
Server	Coordinates system functionality	Section 3.2.
Execution Engine	Enables serverless workflow execution	Section 3.3
Client	Interacts with server and users requests	Section 3.4

Table 1: Laminar Core Elements

Depicted in Figure 3, Laminar architecture distinctly separates the client, server, and execution functionalities. In this serverless architecture, the back-end responds to client-triggered events (captured by the front-end) by executing code within an isolated environment, like a remote computing infrastructure (e.g. Cloud system or HPC cluster). This design ensures scalability and ephemerality, with the environment being dismantled upon completion to optimize resource efficiency.

3.1 Registry

The *Registry*, hosted remotely on the web-based service ⁶, serves as a central repository housing details about users, Processing Elements (PEs), and dispel4py workflows. It includes descriptions and various properties for each entity, as depicted in Figure 4 and Table 2. Users are associated with both PEs and workflows through a one-way many-to-many relationship, ensuring that they are linked to their respective registered PEs or workflows while maintaining data privacy and preventing unauthorized access to information. This design approach promotes the concept of PE and workflow "owners" within the unified registry system, eliminating the need for individual registries. When a user intends to register a PE that is already associated with another user, the system includes the

⁴unixcoder-base

⁵<https://github.com/microsoft/CodeBERT/tree/master/UniXcoder/downstream-tasks/>

⁶<https://www.freesqldatabase.com/>

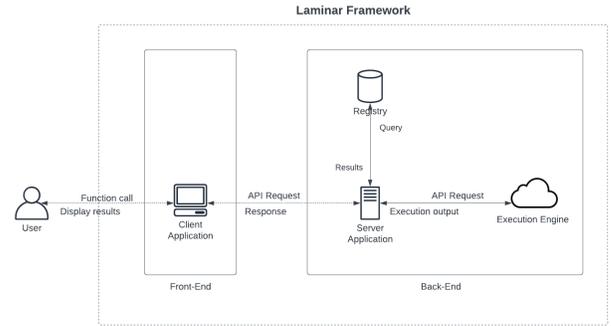


Figure 3: Architectural Overview of Laminar

user as an additional user (or owner) of the PE or workflow, rather than creating a duplicate entry. This process ensures that users can efficiently access their registered PEs or workflows without redundant data entries.

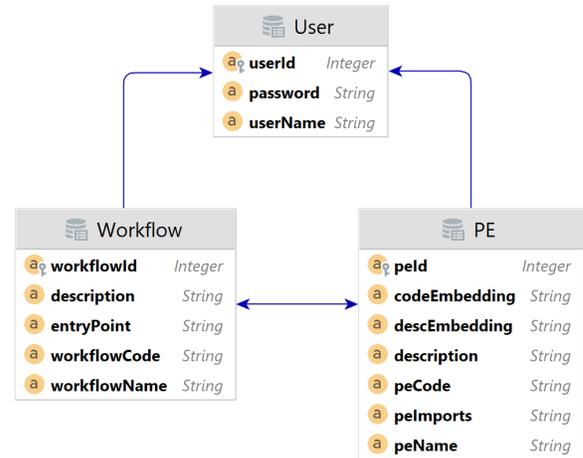


Figure 4: Registry Schema

The relationship between PEs and workflows is established as a two-way many-to-many association. This design allows a PE to be associated with multiple workflows, and conversely, a workflow can consist of numerous PEs. This structure proves beneficial in terms of data duplication and enhanced querying capabilities. For instance, users often require the ability to identify all PEs belonging to a specific workflow. With the current setup, retrieving all PEs associated with a workflow becomes straightforward, simplifying the querying process. This efficient relationship design streamlines user interactions with the system's API, as data pertaining to a workflow can be accessed without the need for additional work.

3.1.1 PE Summarization and Embeddings. To enable seamless semantic code search and code completion, the *Registry* incorporates PE summarizations. These summarizations are generated by the *Client* (refer to Section 4.2) when users do not provide a description

Entity	Properties
PE	<p>peId: A unique identifier for a PE entry</p> <p>codeEmbedding: The code embedding of a PE</p> <p>descEmbedding: The description embedding of a PE</p> <p>description: A field indicating the PE's functionality, provided by the user or automatically summarized</p> <p>peCode: The serialized code for a PE</p> <p>peImports: Import dependencies for a PE</p> <p>peName: The class name of a PE</p>
Workflow	<p>workflowId: A unique identifier for a Workflow entry</p> <p>description: An optional field describing the functionality of a Workflow</p> <p>entryPoint: A unique name identifier for a Workflow</p> <p>workflowCode: The serialized code for a Workflow</p> <p>workflowName: The class name for a Workflow</p>
User	<p>userId: A unique identifier for a user entry</p> <p>password: User password</p> <p>userName: A unique name identifier for a user</p>

Table 2: Registry Entities and Properties

while registering a PE. As introduced in Table 2, the PE **description** property can include either the user-provided description or an automatically generated summary using the *CodeT5* Language model. These PE summaries and descriptions form the cornerstone of our semantic code search functionality.

Additionally, the *Registry* includes storage for PE code embeddings (**codeEmbedding** property) and PE description embeddings (**descEmbedding** property). These embeddings are generated by the *Client* (as discussed in Sections 4.3 and 4.2) during PE registration, contributing to enhanced performance outcomes. Storing these embeddings allows us to perform efficient semantic code searches and code completions without the need to re-calculate them every time a user initiates a search. This re-use of embeddings significantly enhances the responsiveness and performance of our system, ensuring a seamless user experience while navigating and querying the *Registry*.

3.2 Server

The *Server* is a foundational element within the Laminar architecture, housing the core domain logic. Organized using the layered design pattern, it consists of distinct tiers: *Controller*, *Service*, *Model*, and *Data Access Object (DAO)* layers. These layers, pivotal to the architecture, collectively empower the *Server* with modular and structured functionality.

3.2.1 Controller Layer. At the top, the *Controller layer* acts as an interface, handling client requests and orchestrating responses using JSON for data exchange. The *Laminar API* in this layer serves as the conduit for client-server interaction. We have different controllers for different parts of the system (e.g. PE, Workflow, Execution, etc.) Table 3 enumerates the API endpoints for each controller.

3.2.2 Service Layer. Beneath, the *service layer* houses the application's business logic. Here, data from the controller is processed, detached from specific data storage details, and reliant on the DAO layer for data access.

3.2.3 Data Access Object (DAO) layer. The *DAO layer* interacts with the data store, executing essential CRUD operations—create, read, update, delete—ensuring seamless data management.

Controller	Endpoints
PE	<p>POST /registry/{user}/pe/add</p> <p>GET /registry/{user}/pe/all</p> <p>GET /registry/{user}/pe/id/{id}</p> <p>GET /registry/{user}/pe/name/{name}</p> <p>DELETE /registry/{user}/pe/remove/id/{id}</p> <p>DELETE /registry/{user}/pe/remove/name/{name}</p>
Workflow	<p>POST /registry/{user}/workflow/add</p> <p>GET /registry/{user}/workflow/all</p> <p>GET /registry/{user}/workflow/id/{id}</p> <p>GET /registry/{user}/workflow/name/{name}</p> <p>GET /registry/{user}/workflow/pes/id/{id}</p> <p>GET /registry/{user}/workflow/pes/name/{name}</p> <p>DELETE /registry/{user}/workflow/remove/id/{id}</p> <p>DELETE /registry/{user}/workflow/remove/name/{name}</p> <p>PUT /registry/{user}/workflow/{workflowId}/pe/{peId}</p>
Execution	<p>POST /execution/{user}/run</p>
Registry	<p>GET /registry/{user}/all</p> <p>GET /registry/{user}/search/{search}/type/{type}</p>
User	<p>GET /auth/all</p> <p>POST /auth/login</p> <p>POST /auth/register</p>

Table 3: Laminar API Controllers and Endpoints

3.2.4 Model Layer. Lastly, the *Model layer* introduce an object-oriented representation of system data, including *Registry* entity definitions introduced in Section 3.1.

3.2.5 Error Handling. To further augment server-side operations, tailored error handling has been implemented in the *Server* across layers to address unforeseen and unauthorized behaviors, encompassing scenarios such as invalid login credentials. These meticulously designed exceptions furnish critical insights, including type identification, error codes, failed parameters, and supplementary details, thereby enhancing user comprehension and overall experience. Additionally, conformity to a standardized JSON format streamlines exception response formatting on the client side.

3.3 Execution Engine

The *Execution Engine* is the serverless core of Laminar, enabling remote execution of workflows (or just single PEs) with a single API endpoint: /execution/user/run (see Table 3). This endpoint is the gateway for execution requests, including workflows, PEs, runtime configs, arguments, imports, and mappings.

Within a conda Python environment⁷, the execution engine is furnished with the *dispe14py* library and its essential packages. This guarantees a smooth execution environment, sparing users from remote environment management. An intelligent auto-import mechanism scrutinizes the varied import dependencies of streaming workflows, creating an all-inclusive requirement list transmitted to the *Execution Engine*. It autonomously imports necessary prerequisites, eliminating the need for user installations.

The *Execution Engine* handles additional resources essential for workflow execution, such as data from various sources. Users can compile these resources in a 'resources' directory, managed by the application for smooth serialization and deserialization during execution. An additional enhancement is the autonomous identification of the initial PE within a workflow, crucial for invoking the

⁷A conda environment is a self-contained directory that holds specific software packages and their dependencies, allowing for easy management and isolation of software environments.

mapping function that initiates workflow execution. While conventional `dispel4py` often requires users to specify the initial PE, the *Execution Engine* autonomously analyzes the workflow's structure to identify the suitable starting point, minimizing user involvement in complex tasks, aligned with the *Client*'s design philosophy

Leveraging `dispel4py`, Laminar harnesses enhanced capabilities. The *Execution Engine* enables streamlined parallel execution of streaming workflows through diverse `dispel4py` parallel mappings. Users effortlessly select mappings via the *Client*, enhancing application speed and resource utilization.

In the future we plan to expand Laminar's capabilities by enabling the registration of multiple Execution Engines, a process that currently involves manual intervention. This will significantly improve convenience and management efficiency for users.

3.4 Client

The *Client* component serves as a user-friendly Python application, designed to ensure smooth engagement with the Laminar framework. It adopts a dual-layer structure as is illustrated in Figure 5 comprising the *client layer* and the *web_client layer*, enhancing accessibility and ease of use. As follows both layers are explained.

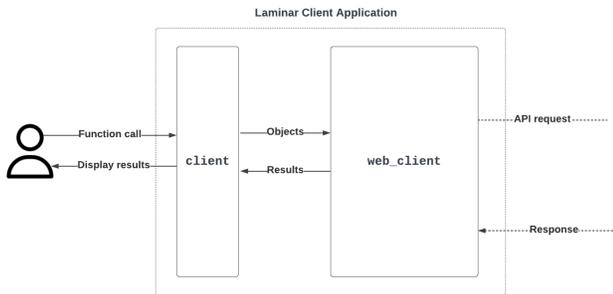


Figure 5: Dual-Layer Client Structure in Laminar Framework

3.4.1 *Client Layer*. This layer offers user-accessible functions for tasks like registering PEs or workflows or executing workflows in a serverless manner. Currently, we support the following functions⁸:

- (1) **register**: User can register with a name and a password.

```
register(self, user_name:str, user_password:str)
```

Example:

```
client.register("zz46", "password")
```

- (2) **login**: Users can login with their details.

```
login(self, user_name:str, user_password:str)
```

Example:

```
client.login("zz46", "password")
```

- (3) **register_PE**: Users can register PEs to store in the *Registry*.

```
register_PE(self, pe:PE_Types, description:str=None)
```

Example:

```
client.register_PE(NumberProducer,
                  "Random numbers producer")
```

- (4) **register_Workflow**: Similarly, workflows can also be registered.

```
register_Workflow(
    self, workflow: WorkflowGraph,
    workflow_name:str,
    description:str=None
)
```

Example:

```
client.register_Workflow(
    graph,
    "isPrime",
    "Workflow that prints random prime numbers"
)
```

- (5) **remove_PE**: PEs in the *Registry* can be deleted if no longer required.

```
remove_PE(self, pe:Union[str,int])
```

Example:

```
client.remove_PE("NumberProducer")
```

- (6) **remove_Workflow**: Users can delete workflows from the *Registry*.

```
remove_Workflow(self, workflow:Union[str,int])
```

Example:

```
client.remove_Workflow("IsPrime")
```

- (7) **get_PE**: registered PEs can be retrieved for creating new workflows.

```
get_PE(self, pe:Union[str,int], describe:bool=False)
```

Example:

```
pe1 = client.get_PE("NumberProducer")
```

- (8) **get_Workflow**: workflows can be retrieved for execution.

```
get_Workflow(
    self,
    workflow:Union[str,int],
    describe:bool = False
)
```

Example:

```
graph = client.get_Workflow("IsPrime")
```

- (9) **get_PEs_By_Workflow**: Users can get a list of PEs for a workflow.

```
get_PEs_By_Workflow(self, workflow:Union[str,int])
```

Example:

```
pes = client.get_PEs_By_Workflow("IsPrime")
```

- (10) **search_Registry**: Search *Registry* for PEs and workflows. **Section 4**

```
search_Registry(
    self,
    search:str,
    search_type:_TYPES = "both",
    query_type:_QUERY_TYPES = "text")
```

Example:

```
results= client.search_Registry("isPrime", "workflow")
```

- (11) **describe**: It provides information on PEs and workflows based on name and description properties.

⁸Full user manual with more examples can be found here <https://tinyurl.com/355zps8p>

```
describe(self, obj:any)
```

Example:

```
client.describe(IsPrime)
```

- (12) **get_Registry**: Retrieves a list of all items stored in the *Registry*.

```
get_Registry(self)
```

Example:

```
registry = client.get_Registry()
```

- (13) **run**: Users can execute workflows at the *Execution Engine*.

```
#process parameter indicates the mapping
# It accepts: SIMPLE, MULTI, MPI, REDIS.

#input indicates the number of iterations
#for which the workflow will be running
run(
    self, workflow:Union[str,int,WorkflowGraph],
    input=None,
    process=_MAPPING_TYPES = "SIMPLE",
    args=None,
    resources=bool=False)
```

Example:

```
#Simple mapping is automatically inferred
# when no mapping is specified.
# Running the workflow for 5 iterations
client.run("IsPrime", input=5)
```

Note that users possess the flexibility to submit workflows using the run function from the client (as demonstrated above), which can involve multiple PEs as showcased in the use cases outlined in Section 5. Alternatively, they have the option to create workflows with a single PE, such as a *ProducerPE* or *GenericPE*, similar to traditional FaaS frameworks. PEs, whether as components of a workflow or individual units run through Laminar, can demonstrate either stateful behavior like the *CountWords* PE in Listing 2, or stateless behavior, as exhibited by the *NumberProducer* in Listing 1. **3.4.2 Web Client Layer.** This layer holds an important role within the Laminar framework. This intermediary layer acts as a conduit for communication, orchestrating the flow of data and interactions between the client-side processing and the server-side execution, as is shown in Figure 5. A significant aspect of this layer is its versatility in handling multiple input types. For instance, the workflow parameter of the run function introduced in 3.4.1 can accept various input formats such as a workflow name (str), a workflow ID (int), or workflow object *WorkflowGraph* for execution - workflow:Union[str,int,WorkflowGraph]. Combining these capabilities enhances user experience and simplifies usage.

A notable challenge that we took in this layer was the code serialization – a crucial step to package Workflows and PEs in a format comprehensible to the execution engine. To tackle this, external packages were utilized for serialization. *cloudpickle* library⁹, chosen after evaluating alternatives like *pickle*¹⁰ and *dill*¹¹, emerged as the preferred option. Its capability to serialize complex Python

```
class CountWords(GenericPE):
    def __init__(self):
        from collections import defaultdict
        GenericPE.__init__(self)
        #Add an input port named "input", from which
        #it will receive tuples with shape (word, 1).
        #Data is group-by (MapReduce)
        #the first element (index 0) of the tuples
        self._add_input("input", grouping=[0])

        #Add an output port named "output"
        self._add_output("output")

        #Initialize a stateful variable
        #to store word counts
        self.count = defaultdict(int)

    def _process(self, inputs):
        import os
        #Extract word and count from the input
        word, count = inputs['input']
        # Update the count for the word
        self.count[word] += count
```

Listing 2: Stateful PE using group-by for word count

objects, including classes and recursive structures, proved essential. Furthermore, *cloudpickle*'s suitability for transmitting code over networks to remote hosts aligned with the project's needs.

Serialization also demanded consideration for storage in the *Registry*. The serialized code, presented as a byte string, needed a suitable format for storage. To ensure portability, a base64 encoding¹² was applied to convert the byte stream into a string format. This serialization approach was leveraged consistently across various code segments required for execution.

This layer also addresses the challenge of handling dependencies. While *cloudpickle* handles import dependencies, an extra library, *findimports*¹³, was used to analyze classes for imports. Users are required to specify necessary imports within PEs to ensure the execution environment has all the required imports for a successful workflow execution. In summary, the *web_client* layer forms the core for client-server interactions and smooths computations throughout Laminar.

4 REGISTRY SEARCH AND EXPLORATION

The *Client* offers a comprehensive set of search and exploration functionalities within the registry, facilitating efficient retrieval and discovery of stored PEs and workflows. This section outlines the different search mechanisms available within Laminar framework.

4.1 Text-Based Search

Text-based searches empower users to quickly locate relevant workflows based on textual information. The *Client* adeptly processes

⁹<https://github.com/cloudpipe/cloudpickle>

¹⁰<https://docs.python.org/3/library/pickle.html>

¹¹<https://pypi.org/project/dill/>

¹²<https://docs.python.org/3/library/base64.html>

¹³<https://pypi.org/project/findimports/>

user-input text queries and effectively matches them with workflows based on their names and descriptions. This functionality includes support for partial matching¹⁴, enabling instances where a user (see Figure 6) queries 'prime' and the system successfully identifies a registered workflow named 'isPrime' (workflow ID 2), thereby enhancing workflow search accuracy and user convenience.

```
client.search_Registry("prime", "workflow")
```

Output:

Searched for "prime"

REGISTRY

Result 1: ID: 2
Workflow Name: isPrime
Description: Workflow that prints random prime numbers

Figure 6: Text-based search for workflows containing the term 'prime' in their names or descriptions. The result is the 'isPrime' workflow with ID '2'.

4.2 Semantic Code Search

In the domain of semantic code search, *LaminarClient* employs a sophisticated process to enable contextually-aware searches for Processing Elements (PEs) based on their user-provided descriptions. This powerful capability enables the discovery of PEs that align with user requirements. If no description is provided during PE registration, the *Client* employs a workaround: it employs code summarization to automatically generate summaries based on the PE's code itself. This automatic summarization utilizes the *codet5-base-multi-sum* model¹⁵, which has been evaluated in comparison to other models for Python code summarization in a recent study [27], resulting in the selection of this model due to its superior performance.

Whether manually input or automatically generated, PE descriptions are transformed into high-dimensional vectors that encapsulate their semantic information, facilitating efficient similarity analysis. This is accomplished using the fine-tuned *unixcoder-code-search* model¹⁶, specialized for code search tasks, as detailed in [8]. Our team fine-tuned this model as an extension of our work [17], adhering to the guidelines outlined in¹⁷, and leveraging the *AdvTest* dataset [20]. Its performance was benchmarked against the *Unixcoder* base model (see Section 6.2.1).

It is noteworthy that the embeddings for PE descriptions are calculated just once, during the registration process. These embeddings are then stored within the *Registry* (refer to Section 3.1) under the **descEmbedding** property.

Upon user query submission, *Laminar* transforms the query using the *unixcoder-code-search* model, calculating cosine similarity against all registered PE description embeddings. This approach,

¹⁴User input text and stored workflow textual information are normalized during a preprocessing step.

¹⁵A fine-tuned variant of the Code T5 model, trained on the CodeSearchNet dataset, available at <https://huggingface.co/Salesforce/codet5-base-multi-sum>

¹⁶<https://huggingface.co/Lazyhope/unixcoder-nine-advtest>

¹⁷<https://github.com/microsoft/CodeBERT/tree/master/UniXcoder/downstream-tasks/>

rooted in the bi-encoders paradigm (explained in Section 2.4), enhances search accuracy and context. The depicted scenario in Figure 7 showcases a user who has registered 22 PEs and five workflows (two detailed in Section 5). In this instance, a semantic search for the text 'A PE that checks if a number is prime' is conducted, resulting in ranked PEs based on similarity scores.

```
client.search_Registry("A PE that checks if  
a number is prime",  
"pe", "text")
```

Output:

Searched for "A PE that check if a number is prime"				
peId	peName	description	similarity	
3	isPrime	check if the input is a prime or not	0.796563	
17	TestProducer	This PE produces a range of numbers	0.502303	
10	InternalExtinction	function to calculate internal extinction from...	0.216504	
19	TestDelayOneInOneOut	This method is called by the PE base class to ...	0.213435	
18	TestOneInOneOut	This PE copies the input to an output	0.186635	

Figure 7: Semantic search of registered PEs. Automated descriptions generated by *Laminar* for PEs with peIDs 3 and 10

This feature enhances the user experience, allowing natural language input to be intelligently matched with relevant PEs. Users later can retrieve (e.g using `get_PE()` client function) a registered PE for creating new workflows.

4.3 Code Completion

Efficient and accurate code completion forms a pivotal cornerstone in the toolkit of not only developers but also scientists and research software engineers. *Laminar* integrates the *ReACC-py-retriever* model (introduced in Section 2.3) to tackle code completion tasks head-on. This functionality empowers users to explore and retrieve relevant Processing Elements (PEs) based on their input. This input can be a partial or complete code query for a specific PE. Our selection of the *ReACC-py-retriever* model is the result of an evaluation against alternative models for the same task (refer to Section 6.2.2 for details).

Much like the process for semantic code search, the *Client* constructs embeddings for each PE's code using the *ReACC-py-retriever* model. These embeddings are subsequently stored in the registry, thoughtfully organized under the **codeEmbedding** property.

When a user seeks to execute a query using a code snippet relevant to a specific PE, the *Client* generates an embedding for the provided query code snippet, employing the same *ReACC-py-retriever* model. The *Laminar* calculates the cosine similarity between the query's embedding and the embeddings of all registered PEs' codes. This comparison leads to the identification of PEs whose code snippets closely resonate with the user's query.

Laminar enhances the coding experience by offering relevant and contextually appropriate PEs. The system accommodates code snippets that may not necessarily warrant completion; they might manifest as fragments of functionality. *Laminar* extracts the most relevant or fitting code snippet to augment the user's input.

```
client.search_Registry("random.randint(1, 1000)",  
"pe", "code")
```

Output:

peId	peName	description	similarity
4	NumberProducer	This function is called to generate a random s...	0.752691
5	PrintPrime	Process the sequence of words in the sequence ...	0.671739
22	TestMultiProducer	Process the sequence of sequence sequen...	0.635483
3	IsPrime	check if the input is a prime or not	0.619670
17	TestProducer	This PE produces a range of numbers	0.552182

Figure 8: Semantic code completion for registered PEs relevant to the code "random.randint(1, 1000)"

Figure 8 shows the process of semantic code completion for registered Processing Elements (PEs) is illustrated. The scenario involves searching for PEs that are relevant to the code snippet `random.randint(1, 1000)`. This code snippet serves as a query for code completion, seeking registered PEs whose code segments are similar to the provided snippet. The results are presented in ascending order of similarity score.

Subsequently, users can retrieve their desired PEs, such as the one with the highest similarity score, to incorporate it into a new workflow or create a new PE by reusing code segments. This feature expedites development, eases cognitive load, and notably minimizes the chances of errors during new PE creation.

5 COMPUTATIONAL SHOWCASES

In this section we are going to introduce two `dispel4py` workflows to showcase Laminar functionality.

5.1 IsPrime workflow

The `IsPrime` workflow (Listing 3) continuously streams random numbers using the 'NumberProducer' PE. It assesses their primality with the 'IsPrime' PE and prints the prime numbers using 'PrintPrime' PE. This exemplifies data flow and processing, connecting PEs to identify and print prime numbers. The workflow is depicted by the green graph in Figure 1.

```

Successfully logged in: root
Successfully executed workflow:
Executing workflow with multi process
before checking data - 982 - is prime or not
before checking data - 801 - is prime or not
before checking data - 797 - is prime or not
IsPrime1 (rank 1): Processed 3 iterations.
before checking data - 390 - is prime or not
before checking data - 179 - is prime or not
IsPrime1 (rank 2): Processed 2 iterations.
the num 797 is prime
the num 179 is prime
PrintPrime2 (rank 4): Processed 1 iteration.
PrintPrime2 (rank 3): Processed 1 iteration.

```

Figure 9: Output sent from the Execution Engine to the Client

When executing the `IsPrime` workflow in Laminar, a mapping strategy must be specified. In this case (shown in Listing 4), we use the `Multi` (process) parallel mapping. The configuration involves five iterations (input) with five processes (num). The first PE generates five random numbers, of which only two are prime in this scenario. Importantly, the `run()` function streamlines not only

```

class NumberProducer(ProducerPE):
    def __init__(self):
        ProducerPE.__init__(self)
    def _process(self):
        # Generate a random number
        result= random.randint(1, 1000)
        # Return the number as the output
        return result

```

```

class IsPrime(IterativePE):
    def __init__(self):
        IterativePE.__init__(self)
    def _process(self, num):
        print("before checking data - %s - \\
            is prime or not" %num)
        #Check if the given input (num) is prime
        if all(num % i != 0 for \\
            i in range(2, num)):
            #Only if the input is prime,
            #the value is return as an output
            return num

```

```

class PrintPrime(ConsumerPE):
    def __init__(self):
        ConsumerPE.__init__(self)
    def _process(self, num):
        # Print the input (num)
        print("the num %s is prime" % num)

```

```

# Create instances of the defined PEs
pe1 = NumberProducer()
pe2 = IsPrime()
pe3 = PrintPrime()

#Create a workflow graph
graph = WorkflowGraph()

#Construct the graph, connecting the PEs
#and their inputs and outputs
graph.connect(pe1, \\
    'output', pe2, 'input')
graph.connect(pe2, \\
    'output', pe3, 'input')

```

Listing 3: IsPrime workflow, which corresponds to the abstract (green) graph shown in Figure 1.

```

client.run(graph, input=5, process=MULTI,
    args={'num':5})

```

Listing 4: Executing IsPrime within Laminar using the specified parameters, which correspond to the concrete (blue) graph shown in Figure 1

serverless workflow execution but also automates the registration of the workflow and its PEs.

Upon receiving the workflow and its configuration, the *Execution Engine* initiates execution with the specified mapping (*Multi* in this case). This configuration automatically adjusts the concrete workflow based on the number of processes chosen, eliminating the need for user intervention. The resulting output (shown in Figure 9) is then forwarded to the *Client* by the *Execution Engine*.

5.2 Astrophysics workflow: Internal Extinction

This workflow has been implemented to calculate the extinction within the galaxies, which is a significant property in astrophysics [5]. This property reflects the dust extinction of the internal galaxies and is used for measuring the optical luminosity¹⁸. This workflow is reusable since it can be regarded as a prior step for other complex tasks which require this property.

In Figure 10, the workflow involves four PEs. The `readRaDec` PE loads coordinate data from an input file, while `getVoTable` downloads a relevant VOTable from the Virtual Observatory website based on those coordinates. The `filterColumns` PE employs the `astropy` library¹⁹ to parse and filter the VOTable data. Lastly, the `internalExt.` PE calculates internal extinction using data from the `Filt` PE. The complete workflow code is available here²⁰. Notably, Laminar detects and installs necessary imports for PEs, ensuring smooth operation within the *Execution Engine*.

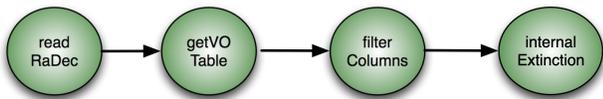


Figure 10: Streaming workflow for calculating the internal extinction of galaxies.

Listing 5, provides an example of how the workflow has been stored with the *Registry*. Following this registration, the workflow becomes accessible for retrieval at any desired moment, as depicted in Listing 6. Ultimately, the workflow can be executed within the serverless *Execution Engine*, as demonstrated in Listing 7. For this instance, we have chosen to execute the workflow using the *Redis* parallel mapping, utilizing ten processes for its execution.

```

client.register_Workflow(
    graph,
    "Astrophysics",
    "A workflow to compute the
    internal extinction of galaxies")
  
```

Listing 5: Registering the Astrophysics workflow

```

workflow = client.get_Workflow("Astrophysics")
  
```

Listing 6: Retrieving the Astrophysics workflow

¹⁸<http://amiga.iaa.es/p/1-homepage.html>

¹⁹Python library for astronomy: <https://www.astropy.org/>

²⁰https://github.com/dispel4pyserverless/dispel4py-client/blob/main/CLIENT_EXAMPLES/AstroPhysics.py

```

client.run(workflow,
    input=[{"input": "resources/coordinates.txt"}],
    process=REDIS,
    args={'num': 10}
    resources=True)
  
```

Listing 7: Executing the Internal Extinction workflow within Laminar, using Redis mapping and 10 processes.

As mentioned in Section 3.3, Laminar supports additional resources that workflows may need. This is exemplified in the *Internal Extinction* workflow (Listing 7), where the required input file, containing relevant coordinate data (e.g., `coordinates.txt`), is accessed from the resources' directory. Through a sequence of copying, serialization, and deserialization steps, the file becomes accessible to the *Execution Engine* for use during workflow execution.

6 EVALUATION

6.1 Laminar Performance

To evaluate the performance of Laminar, a latency analysis was conducted, focusing on both the *Simple* (sequential) and *Multi* mappings (using 5 processes) of the *Internal Extinction* workflow introduced in Section 5.2. This evaluation aimed to compare the execution times between the original `dispel4py` and Laminar, utilizing both local and remote *Execution Engines*. Notably, Laminar accommodates both local and remote *Execution Engines* with some manual adjustments, although in the future we will allow users to register multiple *Execution Engines* from the *Client*.

For the remote deployment scenario, we encapsulated the *Execution Engine* within a Docker image²¹, tailored for easy deployment on diverse cloud platforms. In this experiment, we used the Azure Container Registry²² to host the image, executing it within an Azure web app via Azure App Services²³. This setup provides scalability and adaptability, creating an ideal environment for testing the Laminar framework.

We conducted latency tests for Laminar and compared them to regular `dispel4py` execution. The average execution time for each method and mapping was recorded in seconds. Notably, these tests involved direct execution without workflow registration. However, it is important to note that for both local and remote Laminar experiments, the *Registry* is hosted remotely on the web-based service (see Section 3.1).

Property	Local Ex. Engine	Remote Ex. Engine
OS	Ubuntu 20.04 LTS	Unix 5.15.116.1
Kernel	5.10.16.3-microsoft-standard-WSL2	N/A
CPU	Intel(R) Core(TM) i5-7200U CPU @ 2.50GHz	Intel(R) Xeon(R) CPU E5-2673 v4 @ 2.30GHz
Memory	6.1Gi	1.6Gi

Table 4: Execution Engines Configuration

²¹<https://hub.docker.com/r/zz46/execution>

²²<https://azure.microsoft.com/en-gb/products/container-registry>

²³<https://executionengined4py.azurewebsites.net/run>

Execution Method	Mapping	
	Simple	Multi
<i>original dispel4py</i>	642 sec.	7.32 sec.
<i>Local Execution</i> (with Laminar)	928.2 sec.	11.31 sec.
<i>Remote Execution</i> (with Laminar)	1002 sec.	12.94 sec.

Table 5: Execution times of the *Internal Extinction*

The results in Table 5 reveal a noticeable latency between the original *dispel4py* execution (performed within the same computing environment as the one used for Laminar’s Local Execution Engine) and the utilization of the Laminar framework. Laminar engages in substantial processing, encompassing the retrieval of workflows and their associated PEs from the *Registry*, their preparation for external execution, and an in-depth analysis to ensure dependency management. Additionally, the framework automatically installs the requisite Python libraries within the (local/remote) *Execution Engine*. Furthermore, the request traverses an additional layer at the server, contributing to the observed latency.

Conversely, examining the latency between local and remote *Execution Engines* reveals no substantial increase. This suggests that the framework can harness cloud technology to enhance its performance and usability on a broader scale, without introducing significant performance concerns.

6.2 Deep learning Models

6.2.1 Semantic Code Search. As discussed in Section 2.4, our choice of the fine-tuned *unixcoder-code-search* model for generating embeddings plays a fundamental role in both PE description and user text input for semantic code searches. This decision stems from the model’s proven adeptness in capturing intricate code semantics.

To reinforce our selection, in this work we conducted additional evaluations comparing the fine-tuned *unixcoder-code-search* model with its base version²⁴. Using the Mean Reciprocal Rank (MRR) metric, consistent with prior studies[8], we comprehensively assessed the model’s performance. The detailed results of these evaluations are presented in Table 6, providing further evidence of the model’s effectiveness in enhancing semantic code search capabilities for natural language queries, specifically in the context of zero-shot²⁵ *text-to-code* search.

Model	Zero-shot Code Search	
	CosQA	CSN
	MRR	
<i>unixcoder-base</i>	43.1	44.7
<i>unixcoder-code-search</i>	58.8	72.2

Table 6: Results on zero-shot text-to-code search

The evaluation employed two distinct datasets. CoSQA (Code Search and Question Answering)[11] consists of 20,604 labeled pairs of natural language queries and codes, annotated by at least 3 human annotators. The CSN dataset[10] is derived from the broader CodeSearchNet dataset, with curated filtering of low-quality queries.

²⁴<https://huggingface.co/microsoft/unixcoder-base>²⁵Zero-shot refers to the capability of a model to perform a search tasks without being explicitly trained on the specific queries or examples involved in the search.

The MRR values for both models validate the superior performance of the fine-tuned *unixcoder-code-search* model.

6.2.2 Code Completion. The zero-shot clone detection evaluation by ReACC[19] involved assessing a model’s ability to retrieve similar code segments from a dataset using partial queries, yielding positive outcomes. In this work we have explored different large models for code completion (as discussed in Section 4.3), expanding the assessment to cover a variety of candidates, including an *Unixcoder* model fine-tuned by our team for clone detection (*unixcoder-clone-detection*)²⁶, the fine-tuned *unix-code-search* model, and two state-of-the-art text embedding models: *BAAI/bge-large-en*²⁷ and *thenlper/gte-large*²⁸. Additionally, we evaluated the CodeBERT [3] and GraphCodeBERT [9] models, tailored specifically for source code analysis and comprehension.

Model	MAP@100	Precision at 1
<i>CodeBERT</i>	1.47	4.75
<i>GraphCodeBERT</i>	5.31	15.68
<i>ReACC-retriever-py</i>	9.60	27.04
<i>thenlper/gte-large</i>	1.9	7
<i>BAAI/bge-large-en</i>	8.17	20
<i>unixcoder-clone-detection</i>	10.4	17
<i>unixcoder-code-search</i>	8.53	22.84

Table 7: Zero-shot clone detection evaluation results

For this evaluation, we utilized the CodeNet Python dataset[21], comprising around 14 million code samples as solutions to diverse coding problems. Our analysis employed two key metrics: MAP@100 (Mean Average Precision at 100) and Precision at 1. MAP@100 computes the average precision of the top 100 retrieved items for each query and then calculates the mean across all queries. In contrast, Precision at 1 gauges the model’s accuracy in retrieving the most relevant item as the **top recommendation**. The detailed results of this evaluation are presented in Table 7. Since our primary interest lies in code completion, we place greater importance on the values obtained for Precision at 1. This metric directly signifies the model’s proficiency in retrieving the most similar code from the dataset. As a result, we have chosen the *ReACC-retriever-py* model for code retrieval due to its robust Precision performance.

7 RELATED WORK

Several frameworks in the field of Function-as-a-Service (FaaS) and serverless computing have paved the way for advancements similar to Laminar. We discuss some of the key related frameworks below:

- FuncX [16] specializes in stateless function executions within distributed computing environments. However, it lacks the ability to handle streaming data and support stateful computations, distinguishing it from Laminar’s stream-based processing capabilities.
- PyWren [13] is an open-source FaaS platform for distributing Python functions over the cloud. Like FuncX, it supports stateless function execution and distributed computing, but

²⁶<https://huggingface.co/Lazyhope/unixcoder-clone-detection>²⁷<https://huggingface.co/BAAI/bge-large-en>²⁸<https://huggingface.co/thenlper/gte-large>

lacks integrated streaming data management and deep learning code search features proposed in our work.

- Apache OpenWhisk²⁹ is an open-source FaaS platform known for its flexibility and scalability in event-driven function execution. However, it does not inherently cater to the unique demands of streaming data processing and deep learning code search as does Laminar.
- Apache Flink [14] excels in distributed stream processing with low latency and high throughput, but lacks Laminar’s integrated serverless approach.
- OpenFaaS³⁰ deploys event-driven functions using Docker and Kubernetes, yet lacks Laminar’s comprehensive stream processing and deep learning code search capabilities.

8 CONCLUSIONS AND FUTURE WORK

We present Laminar, a novel serverless stream-based framework with integrated deep learning search features. Unlike conventional serverless platforms, Laminar efficiently manages streaming data, accommodates stateless and stateful computations, and automatically parallelizes streaming applications across various enactment engines. Leveraging the `dispel4py` library, Laminar surpasses existing frameworks by providing a distinct approach to streaming workflow and PE management, complete with automatic library detection and installation. This relieves users from manual installation and management of the execution engine’s dependencies.

Moreover, Laminar integrates advanced large language models, introducing powerful deep learning code search and completion capabilities. This integration enables sophisticated PE code search through both code and natural language, facilitating seamless exploration and utilization of an extensive repository of PE functionalities. The framework also incorporates code summarization, automating PE functionality summaries and enhancing semantic search. Moving forward, we aim to integrate Laminar with popular cloud providers, enable multiple *Execution Engine* registration, explore dynamic resource provisioning and container management strategies and enhance deep learning search for workflows.

REFERENCES

- [1] 2023. A comprehensive review of State-of-The-Art methods for Java code generation from Natural Language Text. *Natural Language Processing Journal* 3 (2023), 100013. <https://doi.org/10.1016/j.nlp.2023.100013>
- [2] Dirk Eddelbuettel. 2022. A Brief Introduction to Redis. arXiv:2203.06559 [stat.CO]
- [3] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiao Cheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. arXiv:2002.08155 [cs.CL]
- [4] Rosa Filgueira, Amrey Krause, Malcolm Atkinson, Iraklis Klampanos, and Alexander Moreno. 2016. `dispel4py`: A Python Framework for Data-Intensive Scientific Computing. *International Journal of High Performance Computing Applications (IJHPCA)* (2016).
- [5] Rosa Filgueira, Amrey Krause, Alessandro Spinuso, Iraklis Klampanos, Peter Danecsek, and Malcolm Atkinson. 2015. `Dispel4py`: An Open-Source Python library for Data-Intensive Seismology. *EGUGA* (2015), 6790.
- [6] Message P Forum. 1994. *MPI: A Message-Passing Interface Standard*. Technical Report. USA.
- [7] Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. 2018. Deep Code Search. In *Proceedings of the 40th International Conference on Software Engineering (Gothenburg, Sweden) (ICSE '18)*. Association for Computing Machinery, New York, NY, USA, 933–944. <https://doi.org/10.1145/3180155.3180167>
- [8] Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. 2022. UniXcoder: Unified Cross-Modal Pre-training for Code Representation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22–27, 2022*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, 7212–7225. <https://doi.org/10.18653/v1/2022.acl-long.499>
- [9] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, Michele Tufano, Shao Kun Deng, Colin Clement, Dawn Drain, Neel Sundaresan, Jian Yin, Daxin Jiang, and Ming Zhou. 2020. GraphCodeBERT: Pre-training Code Representations with Data Flow. <https://doi.org/10.48550/ARXIV.2009.08366>
- [10] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, Michele Tufano, Shao Kun Deng, Colin Clement, Dawn Drain, Neel Sundaresan, Jian Yin, Daxin Jiang, and Ming Zhou. 2021. GraphCodeBERT: Pre-training Code Representations with Data Flow. arXiv:2009.08366 [cs.SE]
- [11] Junjie Huang, Duyu Tang, Linjun Shou, Ming Gong, Ke Xu, Daxin Jiang, Ming Zhou, and Nan Duan. 2021. CoSQA: 20, 000+ Web Queries for Code Search and Question Answering. *CoRR abs/2105.13239* (2021). arXiv:2105.13239 <https://arxiv.org/abs/2105.13239>
- [12] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. CodeSearchNet Challenge: Evaluating the State of Semantic Code Search. <https://doi.org/10.48550/ARXIV.1909.09436>
- [13] Eric Jonas, Shivaram Venkataraman, Ion Stoica, and Benjamin Recht. 2017. Occupy the Cloud: Distributed Computing for the 99%. *CoRR abs/1702.04024* (2017). arXiv:1702.04024 <https://arxiv.org/abs/1702.04024>
- [14] P Carbone Asterios Katsifodimos, S Ewen Volker Markl, and S Haridi Kostas Tzoumas. 2015. Apache FlinkTM: Stream and batch processing in a single engine. *Bull. IEEE Comput. Soc. Tech. Comm. Data Eng* 36, 4 (2015).
- [15] Manoj Kumar. 2019. Serverless architectures review, future trend and the solutions to open problems. *American Journal of Software Engineering* 6, 1 (2019), 1–10.
- [16] Zhuozhao Li, Ryan Chard, Yadu Babuji, Ben Galewsky, Tyler J. Skluzacek, Kirill Nagaitsev, Anna Woodard, Ben Blaiszik, Josh Bryan, Daniel S. Katz, Ian Foster, and Kyle Chard. 2022. funcX: Federated Function as a Service for Science. *IEEE Transactions on Parallel and Distributed Systems* 33, 12 (dec 2022), 4948–4963. <https://doi.org/10.1109/tpds.2022.3208767>
- [17] Zihao Li and Rosa Filgueira. 2023. Mapping the repository landscape: harnessing similarity with RepoSim and RepoSnipy. In *2023 IEEE 19th International Conference on e-Science (e-Science)*. IEEE. <https://www.escience-conference.org/2023/> 19th IEEE International Conference on eScience, eScience ; Conference date: 09-10-2023 Through 13-10-2023.
- [18] Liang Liang, Rosa Filgueira, Yan Yan, and Thomas Heinis. 2022. Scalable adaptive optimizations for stream-based workflows in multi-HPC-clusters and cloud infrastructures. *Future Generation Computer Systems* 128 (2022), 102–116. <https://doi.org/10.1016/j.future.2021.09.036>
- [19] Shuai Lu, Nan Duan, Hojae Han, Daya Guo, Seung won Hwang, and Alexey Svyatkovskiy. 2022. ReACC: A Retrieval-Augmented Code Completion Framework. arXiv:2203.07722 [cs.SE]
- [20] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin B. Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. 2021. CodeXGLUE: A Machine Learning Benchmark Dataset for Code Understanding and Generation. *CoRR abs/2102.04664* (2021).
- [21] Ruchir Puri, David S. Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi, Vladimir Zolotov, Julian Dolby, Jie Chen, Mihir Choudhury, Lindsey Decker, Veronika Thost, Luca Buratti, Saurabh Pujar, Shyam Ramji, Ulrich Finkler, Susan Malaika, and Frederick Reiss. 2021. CodeNet: A Large-Scale AI for Code Dataset for Learning a Diversity of Coding Tasks. arXiv:2105.12655 [cs.SE]
- [22] Johann Schleier-Smith, Vikram Sreekanti, Anurag Khandelwal, Joao Carreira, Neeraja J. Yadwadkar, Raluca Ada Popa, Joseph E. Gonzalez, Ion Stoica, and David A. Patterson. 2021. What Serverless Computing is and Should Become: The next Phase of Cloud Computing. *Commun. ACM* 64, 5 (apr 2021), 76–84. <https://doi.org/10.1145/3406011>
- [23] Hossein Shafiei, Ahmad Khonsari, and Payam Mousavi. 2022. Serverless computing: a survey of opportunities, challenges, and applications. *Comput. Surveys* 54, 11s (2022), 1–32.
- [24] Shuai Wang, Jinyang Liu, Ye Qiu, Zhiyi Ma, Junfei Liu, and Zhonghai Wu. 2019. Deep learning based code completion models for programming codes. In *Proceedings of the 2019 3rd International Symposium on Computer Science and Intelligent Control*. 1–9.
- [25] Xin Wang, Yasheng Wang, Fei Mi, Pingyi Zhou, Yao Wan, Xiao Liu, Li Li, Hao Wu, Jin Liu, and Xin Jiang. 2021. SynCoBERT: Syntax-Guided Multi-Modal Contrastive Pre-Training for Code Representation. <https://doi.org/10.48550/ARXIV.2108.04556>

²⁹<https://openwhisk.apache.org/documentation.html>

³⁰<https://docs.openfaas.com/>

- [26] Yue Wang, Weishi Wang, Shafiq Joty, and Steven C. H. Hoi. 2021. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation. <https://doi.org/10.48550/ARXIV.2109.00859>
- [27] Christopher Williams and Rosa Filgueira. 2023. RepoGraph: a novel semantic code exploration tool for python repositories based on knowledge graphs and deep learning. In *2023 IEEE 19th International Conference on e-Science (e-Science)*. IEEE. <https://www.escience-conference.org/2023/> 19th IEEE International Conference on eScience, eScience ; Conference date: 09-10-2023 Through 13-10-2023.
- [28] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. HuggingFace's Transformers: State-of-the-art Natural Language Processing. *CoRR* abs/1910.03771 (2019). arXiv:1910.03771
- [29] Chunyan Zhang, Junchao Wang, Qinglei Zhou, Ting Xu, Ke Tang, Hairen Gui, and Fudong Liu. 2022. A Survey of Automatic Source Code Summarization. *Symmetry* 14, 3 (2022). <https://www.mdpi.com/2073-8994/14/3/471>