# W

# A Prediction based Autoscaling in Serverless Computing

Ha-Duong Phung, Younghan Kim School of Electronic Engineering, Soongsil University, Seoul, Korea

> TCSS 562 Team 15: Yuan Huang, Yifan Xie, Alan Liu

#### UNIVERSITY of WASHINGTON



- > Introduction
- > **Describe the Knative and analyze some related forecasting models.**
- > Detail of the proposed method structure
- Presents the results showing how different concurrency limits impact performance and the efficiency of their prediction model Bi-LSTM compared to the Knative default.
- > Conclusion and the future work plans.

#### UNIVERSITY of WASHINGTON



3

# **Introduction:** Paper Overview

#### Why is it a problem?

- 1. Users have no idea about the appropriate parameter values, while those predefined values can strongly impact the performance.
- 2. The negative influence can cause response time increase and throughput reduction if those are overprovisioning or under-provisioning.
- 3. Knative uses a moving average method to calculate the number of pods based on past data but cannot reflect the future workload trend, leading to the delay effect.

#### Why is it a problem the research community is interested?

Due to the fast growth of serverless computing offerings, they focus on maximizing performance with the lowest resource utilization and optimizing response time to satisfy the quality of service (QoS) requirements so that it is content with users' experience while users only must pay the least cost.

### **Introduction:** Paper Overview

#### Address the problem:

Automatically optimizes an effective scaling policy to a specific application.

#### How?

- With the same amount of resources per pod, the purpose focuses on outputting an appropriate parameter level showing the throughput maximum and an acceptable latency to satisfy SLO violation.
- > A new service revision with the new scheduling policy is created, and then apply the forecasting model de to optimize the calculation of the number of pods.

UNIVERSITY of WASHINGTON

# BACKGROUND

- A. Knative Serverless Platform
- B. Forecasting model













# **Proposed Architecture**

- > Two primary states to address above challenges.
  - First State: Parameter values for the primary service function to optimize performance & the latency of service
  - Second State: Improve the calculation of the number of pods to be more adaptive to the changes in request trends by applying the prediction model Bidirectional LSTM(Bi-LSTM)
- > The activator controls these two states, which run separately

UNIVERSITY of WASHINGTON







# **Proposed Architecture**

#### Parameter Values Optimization

- > Run performance tests for varying concurrency levels
- > Collect the output information:
  - maximum throughput, request latency distribution, and corresponding concurrency level for each test iteration
- > Begin at a concurrency limit of 10 and use steps of 10 to update the next limit. Stop when latency is over SLO violation
- Create a new service revision using the output and transfer 100% workload to the latest revision









## **Preliminary Experimental Result**

- The number of pods obtained based on the Bi-LSTM model is closer to the number of ho
- The predicted workload from the Bi-LSTM model is more accurate than the moving average, so the predicted number of pods value can be quickly adaptive to the workload trend



## Conclusion

- > Find the appropriate parameter values for Knative service revision to improve performance and keep acceptable latency
- > Optimize the calculation of the number of pods for Knative based on the prediction model Bi-LSTM
- > The preliminary experiments worked well and show better performance than Knative scheme

UNIVERSITY of WASHINGTON

23

# Strengths

- > Clear logic in paper design
- > The underlying principle and implementation logic of the algorithm and forecasting models are introduced in detail
- > The algorithm has a substantial improvement (over 20%) in performance and the effect is shown in the graphs
- Well-designed and described experimental architecture and controls







