

Not All Explorations Are Equal: Harnessing Heterogeneous Profiling Cost for Efficient MLaaS Training

Jun Yi¹, Chengliang Zhang², Wei Wang², Cheng Li³, Feng Yan¹

¹University of Nevada,
Reno



²Hong Kong University of
Science and Technology



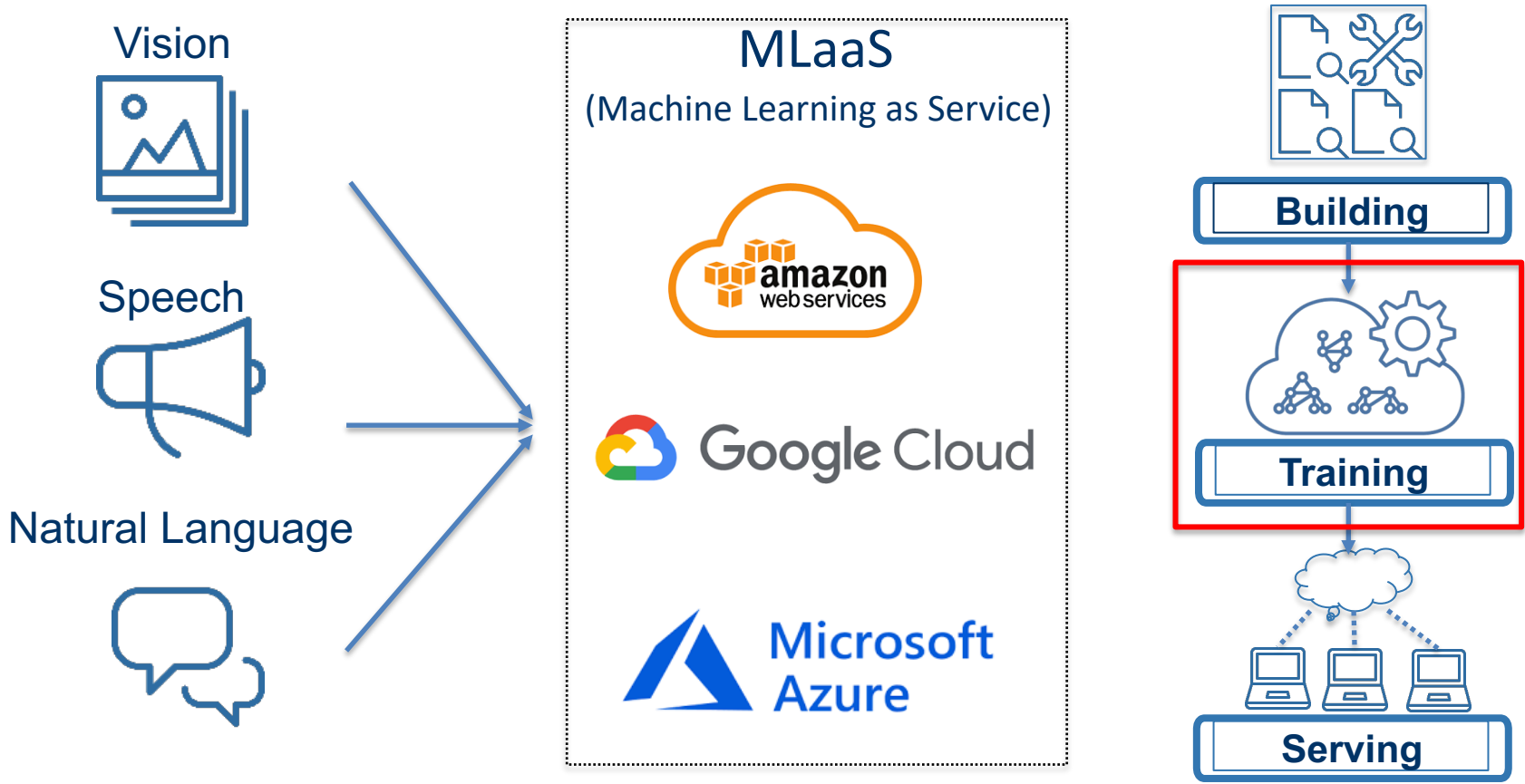
³University of Science and
Technology of China



Resource Acknowledgement



OIT | Cyberinfrastructure



Practical MLaaS training scenarios:

- *Scenario-1*: Training project without time or cost limit
- *Scenario-2*: Training project with time limitation
- *Scenario-3*: Training project with cost limitation

How to deploy MLaaS training jobs in Cloud?

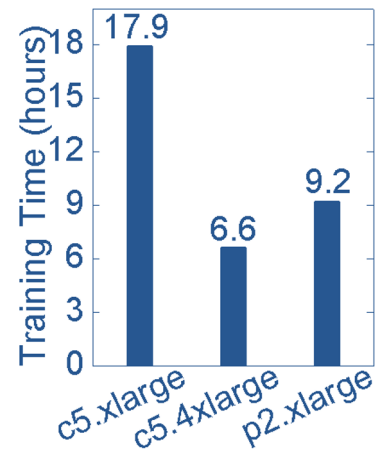
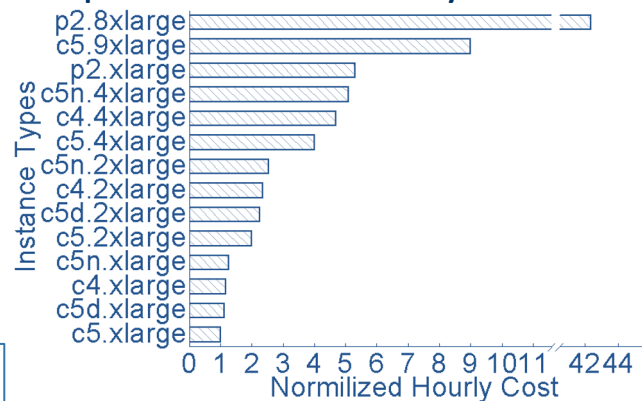
Scale-up (more capable instance)

VS scale-out (more instances)

E.g., use many cheapest instances (40 c5.4xlarge) or a few costly instances (9 p2.xlarge)?

Neither case is optimal (see the right figure)

Up to 42.5X in hourly cost



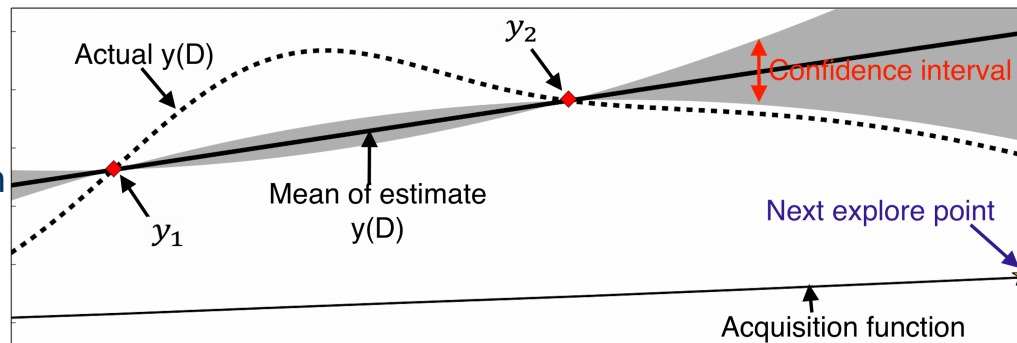
Challenges: large deployment scheme search space (62 scale-up & 50 scale-out->3100 schemes)

Existing Work

<p>Analytical Modeling (assumptions on model/hardware)</p>	<ul style="list-style-type: none"> Limited applicability (fast-evolving ML models) Poor fit for cloud (increasing diversified hardware) 	<ul style="list-style-type: none"> ❖ [SIGKDD '15] <i>Performance modeling and scalability optimization of distributed deep learning systems</i> ❖ [ICLR, '17] <i>Paleo: A performance model for deep neural networks.</i>
<p>Reinforcement Learning</p>	<ul style="list-style-type: none"> Requires extensive training samples and high computing resources 	<ul style="list-style-type: none"> ❖ [Nature '15] Human-level control through deep reinforcement learning.
<p>Pareto-Optimization</p>	<ul style="list-style-type: none"> Falls short in performance 	<ul style="list-style-type: none"> ❖ [CCGRID '17] Predicting cloud performance for hpc applications: A user-oriented approach.
<p>Conventional Bayesian Optimization (BO) (assume uniform profiling cost of every point)</p>	<ul style="list-style-type: none"> Assume uniform exploration cost Lack of ML-specific insights 	<ul style="list-style-type: none"> ❖ [NSDI '17] <i>CherryPick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics.</i> ❖ [ICDCS '18] <i>Arrow: Low-level augmented bayesian optimization for finding the best cloud vm.</i>

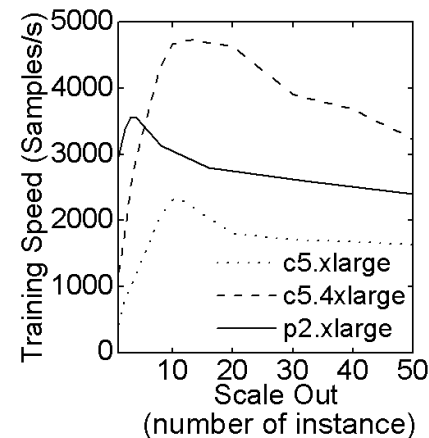
Conventional BO:

- For problems with unknown objective function
- Start with random initial points
- Select next points based on acquisition function
- Acquisition function optimizes expected improvement, probability of improvement, confidence bound, etc.



Key Observations

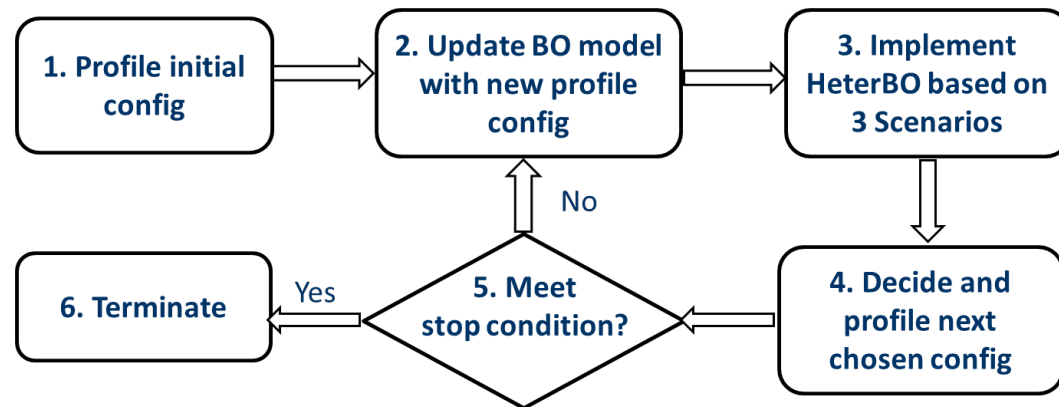
- Heterogenous exploration cost
 - Some schemes (i.e., large scale-out, high-end GPU instance) are more costly to explore than others
- No ML-specific prior is adopted in deployment optimization
 - Speedup trend of scale-out follows a concave-shape curve



Main Idea: Heterogenous cost-aware and ML prior aware BO

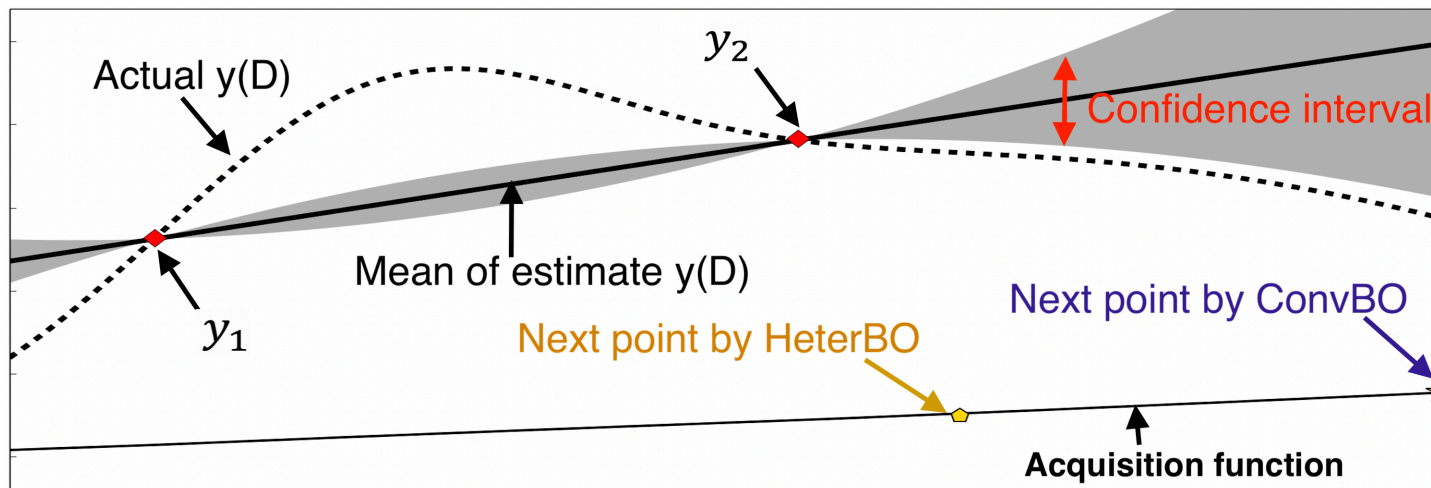
- Problem formulation minimize $T(D)/C(D)$ $T(D)$ - Total Time; $C(D)$ - Total Cost
 subject to $D \in D(m, n)$ $D(m, n)$ - Possible schemes; m - Instance type
 n - Number of selected Instance type

- Search process



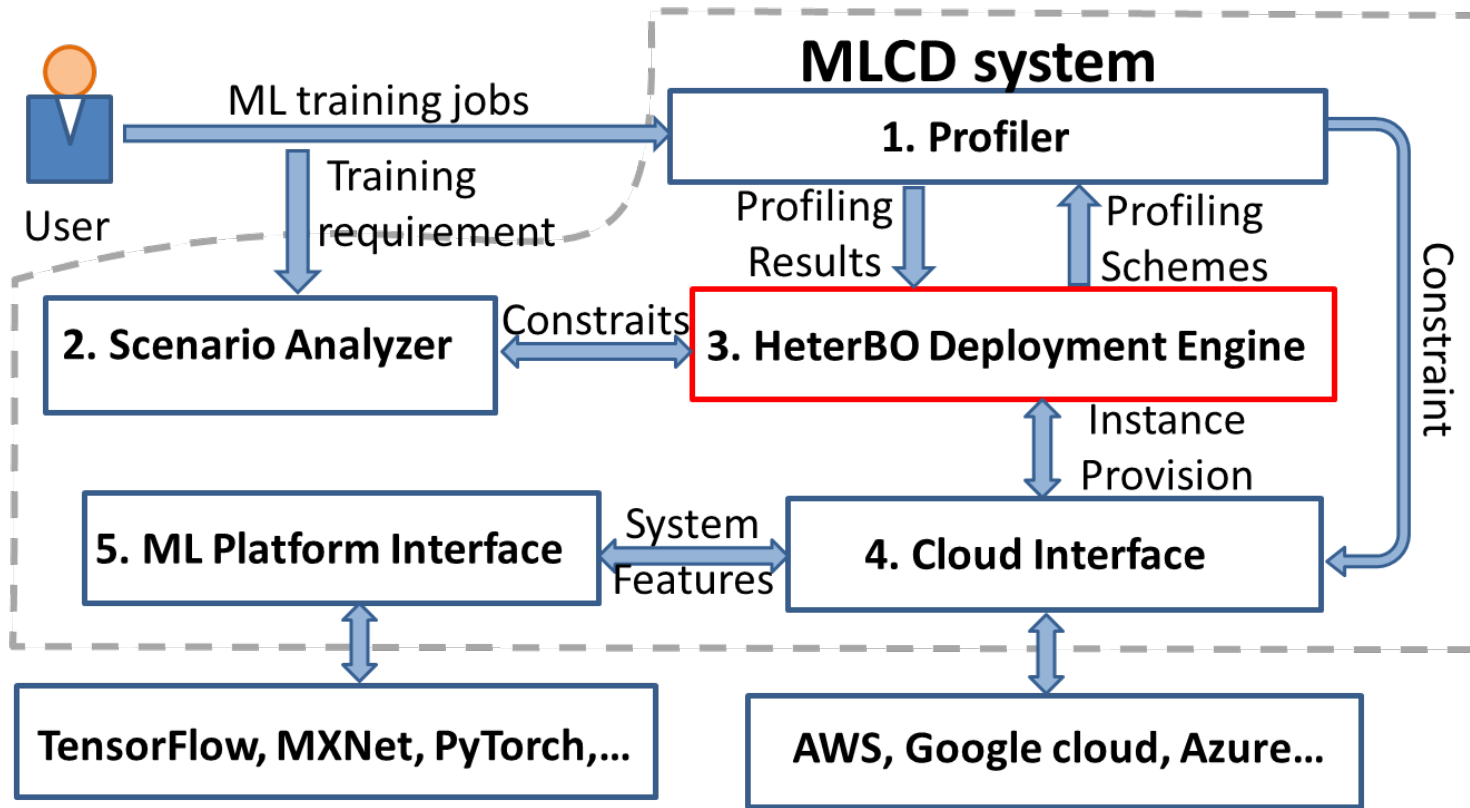
- Key Components

- Prior function: Gaussian Process (flexibility and tractability)
- Acquisition function: EI (Expected Improvements) with constraints (profiling cost) -> $T(\text{rue})EI$
- Heterogeneous search cost aware: avoid randomly jumping into high profile cost regions
- ML-specific aware: detects down slope of the concave-shape -> avoid high overheads



- y_1 and y_2 are profiled points
- Not select the maximum point in acquisition function as next point (i.e., ConvBO)
- HeterBO considers the *user constraints* and *heterogeneous search cost* when selecting next point (35% less profiling cost)

MLaaS training Cloud Deployment system (MLCD):



Testbed

- AWS CPU, GPU instances

ML platforms

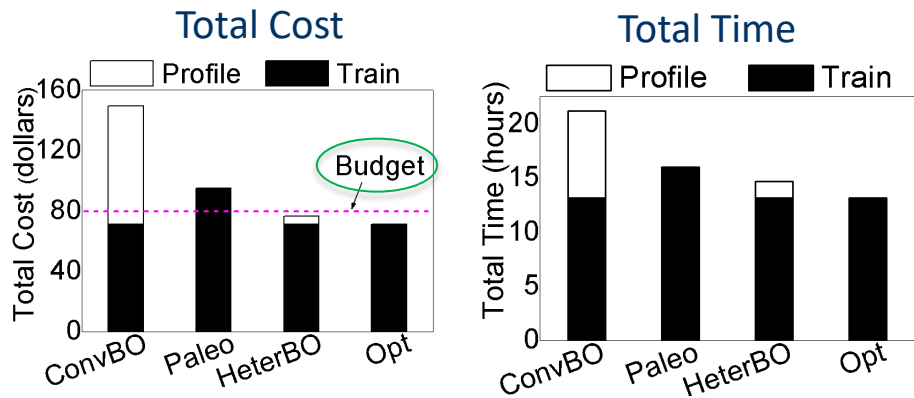
- TensorFlow and MXNet

ML Models

- AlexNet, ResNet, Inception-v3, CharCNN, BERT

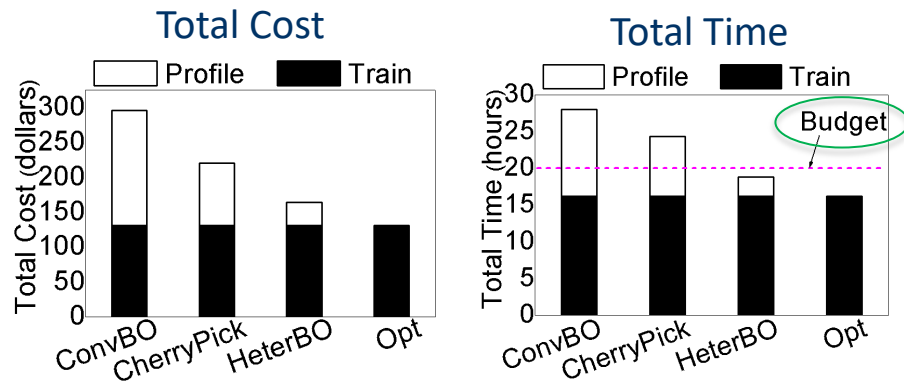
HeterBO vs. Existing Approaches using TensorFlow

Limited monetary budget (\$80) scenario



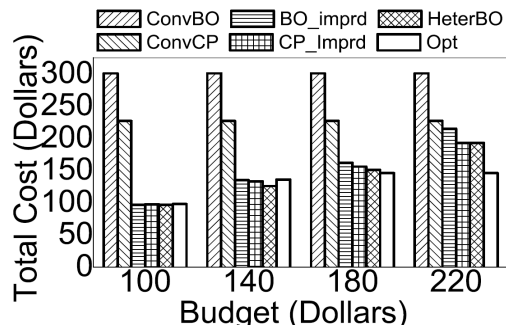
HeterBO costs under budget (ConvBO/Paleo not)
36.4% and **12.5%** better than ConvBO and Paleo
in Total Time

Limited total time (20 hours) scenario



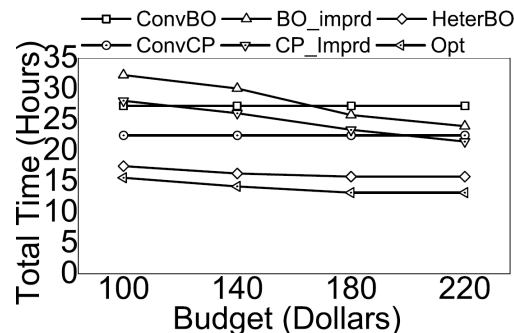
HeterBO finishes on time (ConvBO/CheryPick not)
44.8% and **28.9%** better than ConvBO and CheryPick
in Total Cost

Total cost vs Budget



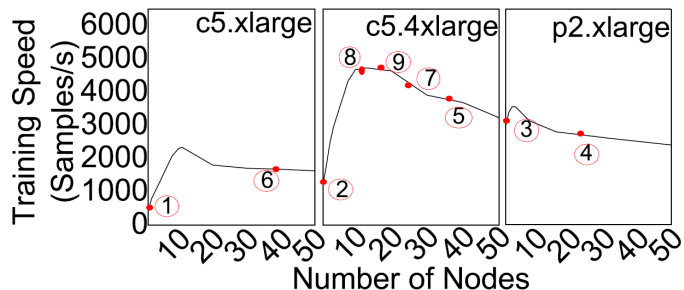
HeterBO outperforms SOTA by up to **3.1x**

Total Time vs Budget



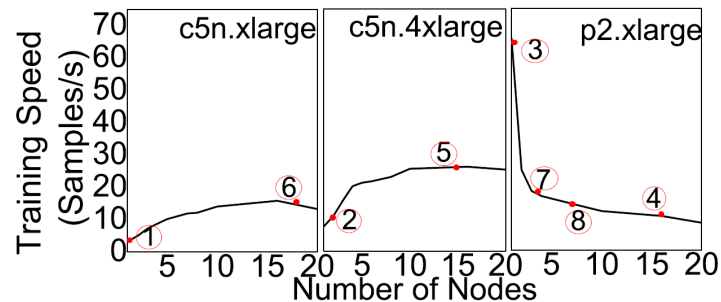
HeterBO outperforms SOTA by up to **2.34x**

Char-RNN using TensorFlow



HeterBO found optimal within budget \$120

BERT using MXNet



HeterBO found optimal within budget \$120

Takeaway:

Not all explorations are equal: heterogeneous exploration cost + machine learning specific prior

→ A fully-automated MLaaS training Cloud Deployment system (**MLCD**) driven by **HeterBO** search method

Jun Yi

junyi@nevada.unr.edu, <https://www.cse.unr.edu/~jyi/>

https://www.youtube.com/channel/UCMgXRQdpjImc5GLkGV0Av8g?view_as=subscriber



OIT | Cyberinfrastructure

Research Grants Council
of Hong Kong
香港 研究資助局

