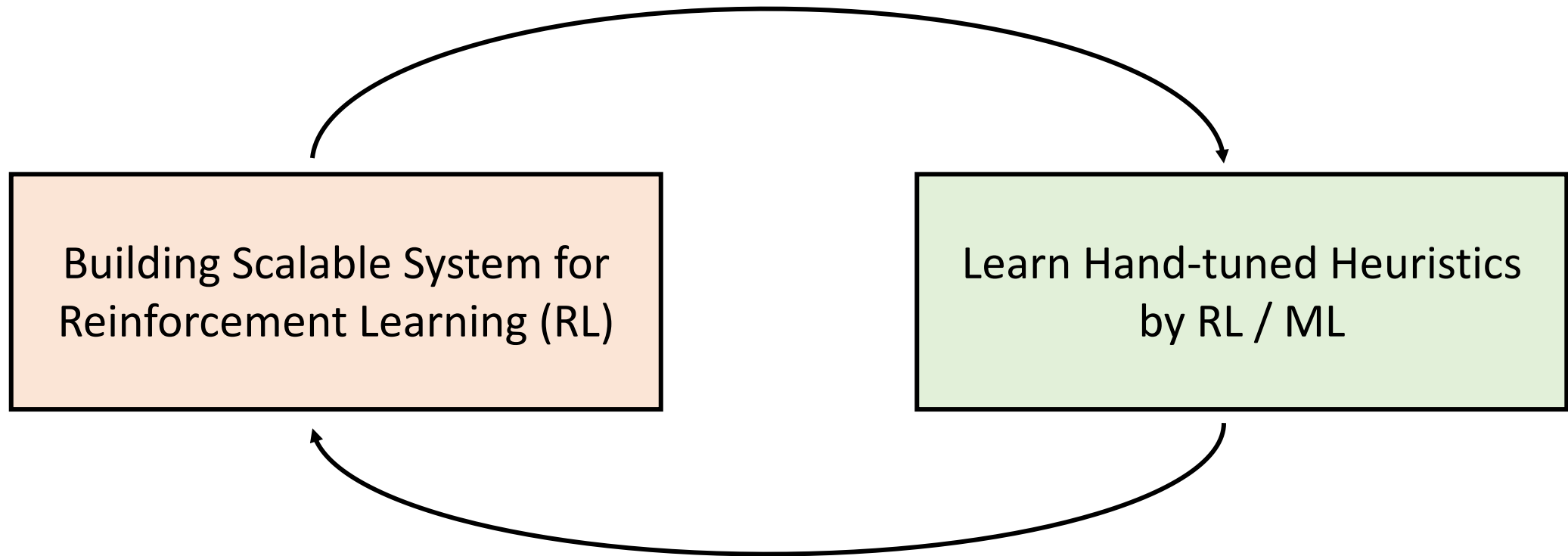


# Building Scalable Systems for Reinforcement Learning and Using Reinforcement Learning for Better Systems

Presented by Yuandong Tian

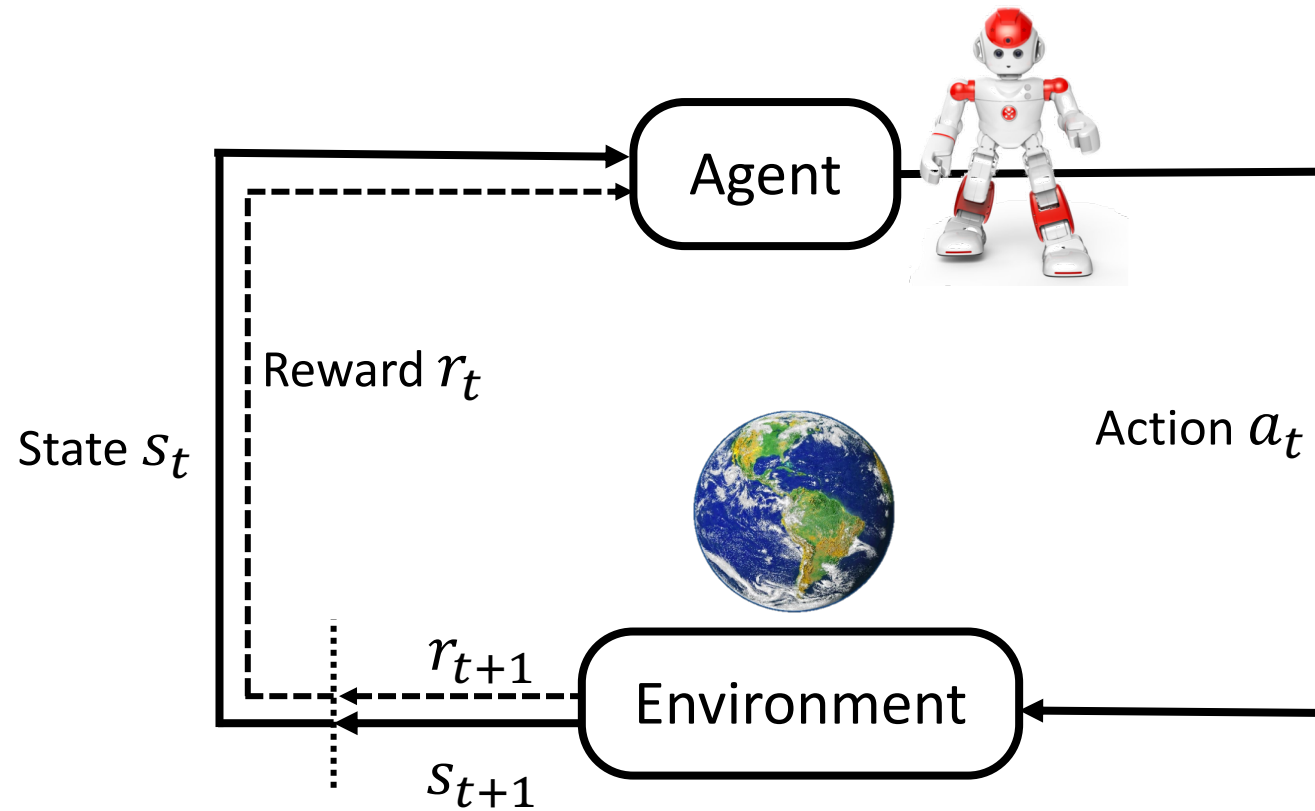
Research Scientist and Manager  
Facebook AI Research

# Overview



# Building Scalable System for RL

# Crash Course of Reinforcement Learning



# Reinforcement Learning works, but expensive

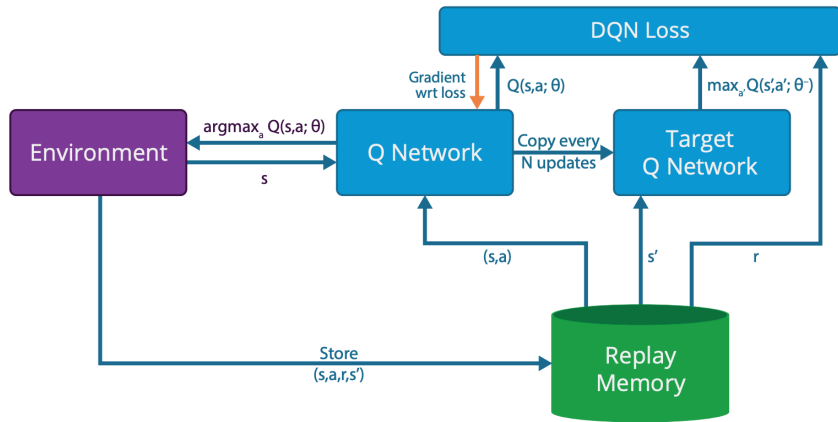


Year	Projects	Human Data	Training Resource	Training time
2016	DeepMind's AlphaGo	Yes	~50 GPUs + ? CPUs	~1 week
2017	DeepMind's AlphaGo Zero (20 blocks)	No	~2000 TPUs	3 days
2017	DeepMind's AlphaZero (20 blocks)	No	~5000 TPUs	8 hours
2018	OpenAI Five	No	128,000 CPUs + 256 GPUs	Several months
2019	DeepMind's AlphaStar	Yes	16,000 CPUs + 3072 TPUv3 cores	44 days

# Challenges in large-scale RL Training System

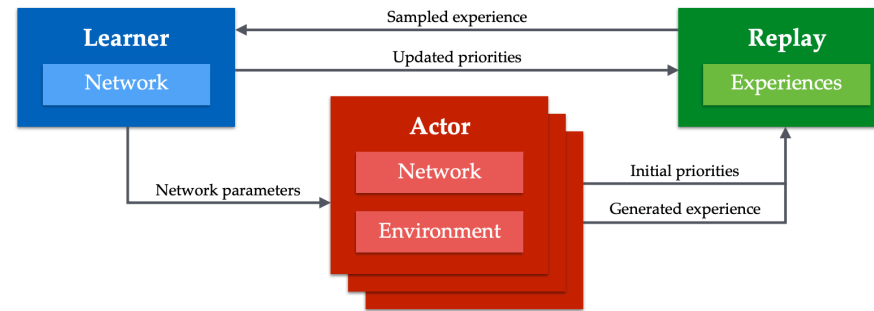
- Trade-offs in a *heterogenous* system
  - **Different kind of objects:** Actor / Environment / Trainer / Replay buffer
  - CPUs / GPUs Allocations
  - Multi-threading versus Multiple Processes, Batching issues
  - Local versus Distributed
  - Synchronization / Asynchronization.
    - On-policy versus off-policy methods
    - Perfect synchronization might NOT give you the best performance
- Mingled Algorithm Design and System Design
  - New System design  $\leftrightarrow$  New RL algorithm

# Distributed System for training RL agent



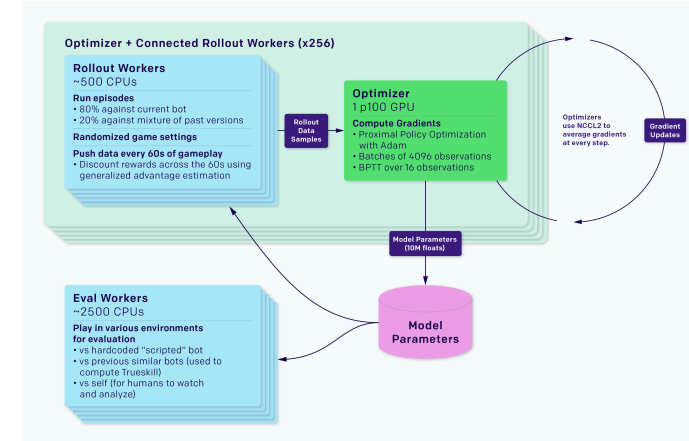
GORILLA

[Massively Parallel Methods for Deep Reinforcement Learning, AAI 2015]



Ape-X / R2D2

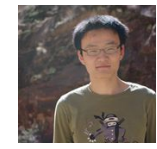
[Distributed Prioritized Experience Replay, Horgan et al, ICLR 2018]  
 [Recurrent Experience Replay in Distributed Reinforcement Learning Kapturowski et al, ICLR 2019]



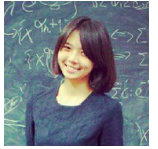
OpenAI Rapid



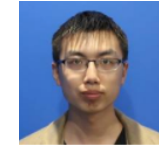
Yuandong Tian



Qucheng Gong



Wendy Shang



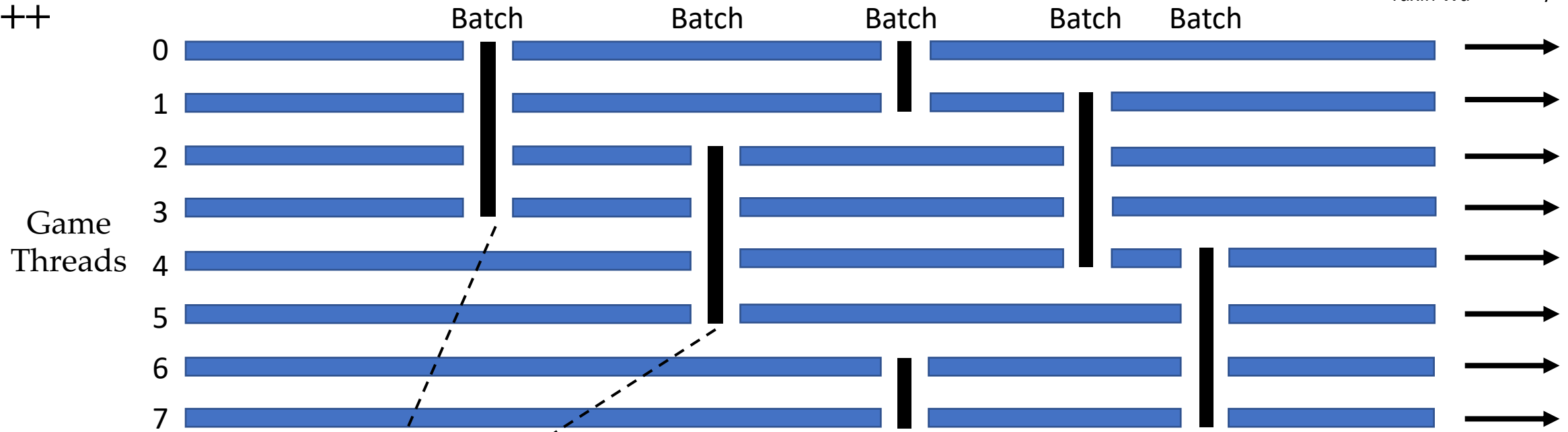
Yuxin Wu



Larry Zitnick

# ELF: RL Framework for Game Research

C++



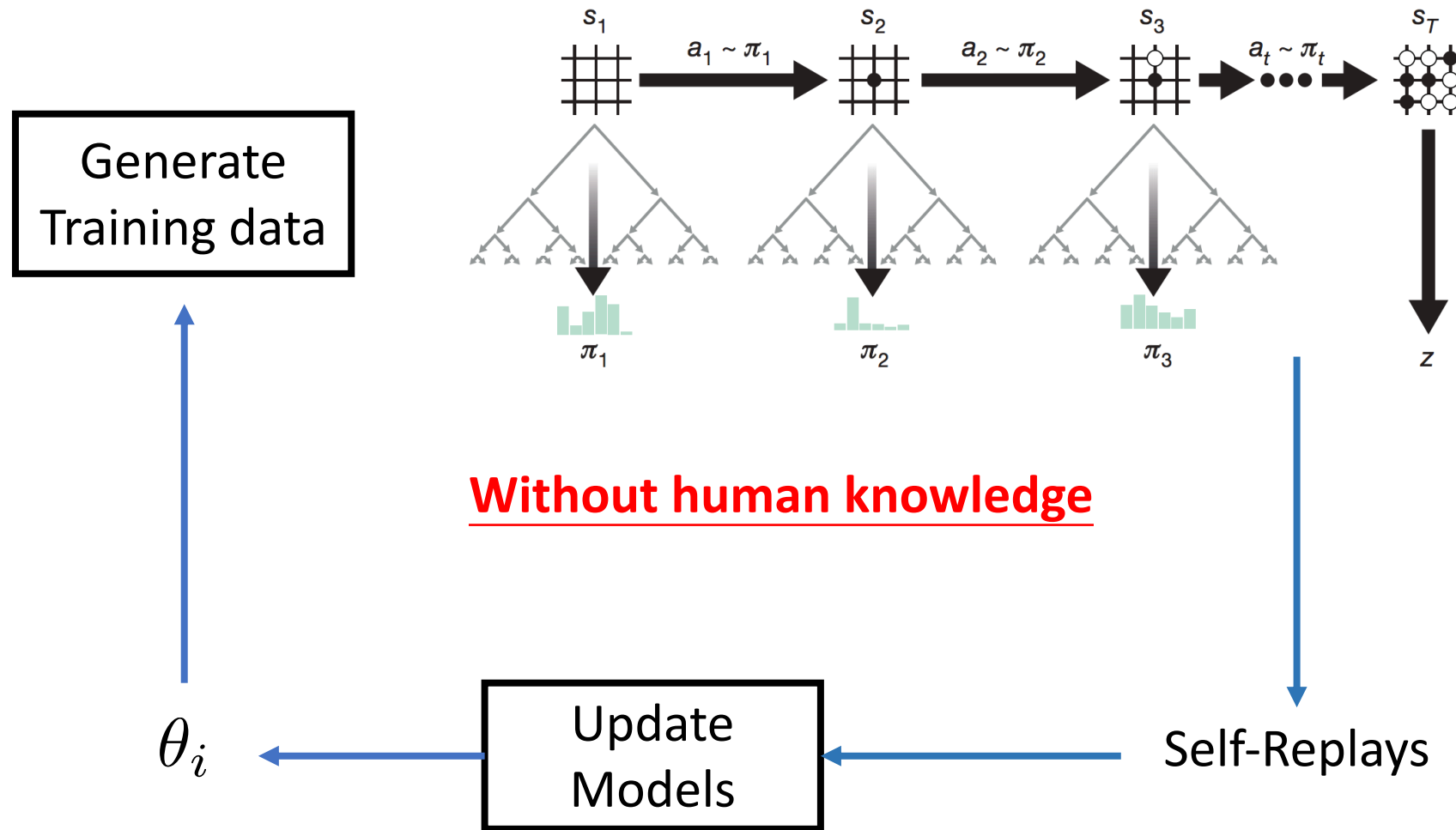
Game  
Threads

Python

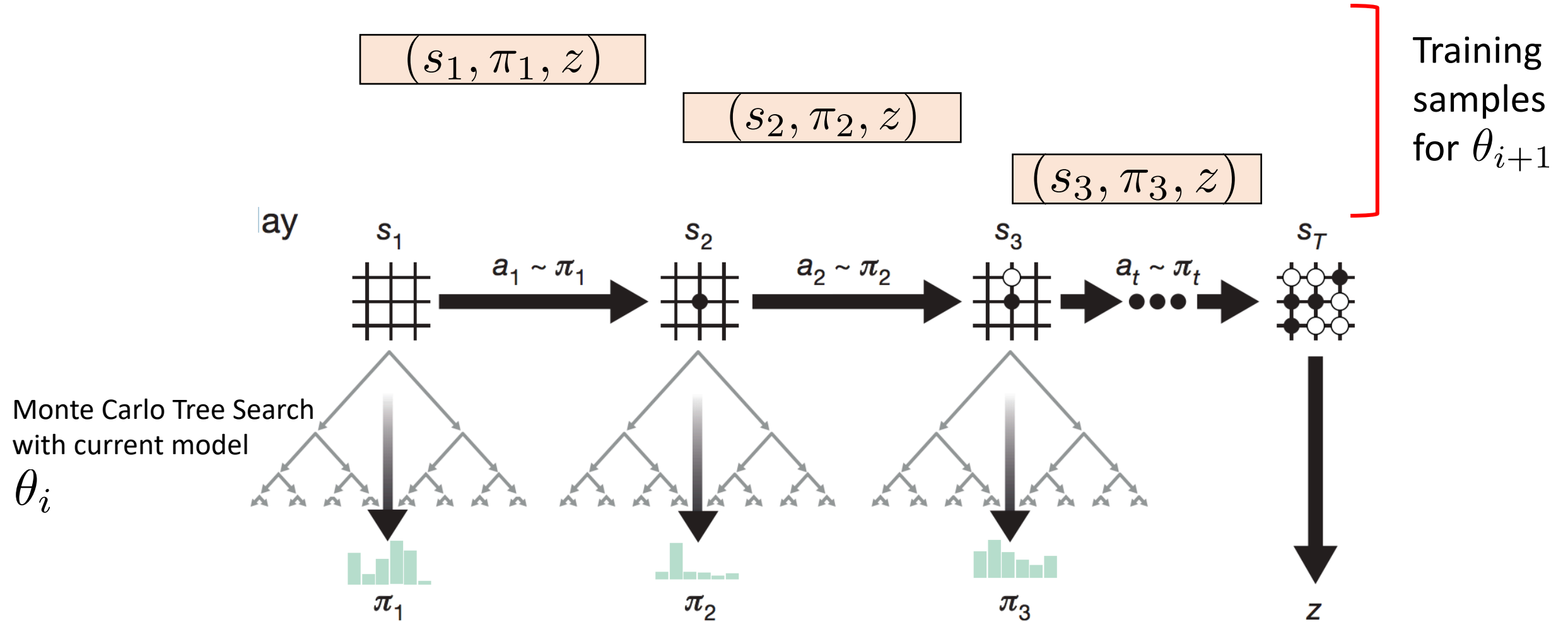
```
while True:
    batched_states = GameContext.Wait()
    replies = model(batched_states)
    GameContext.Steps(replies)
```



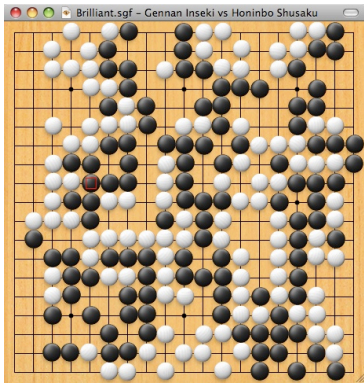
# AlphaGoZero / AlphaZero



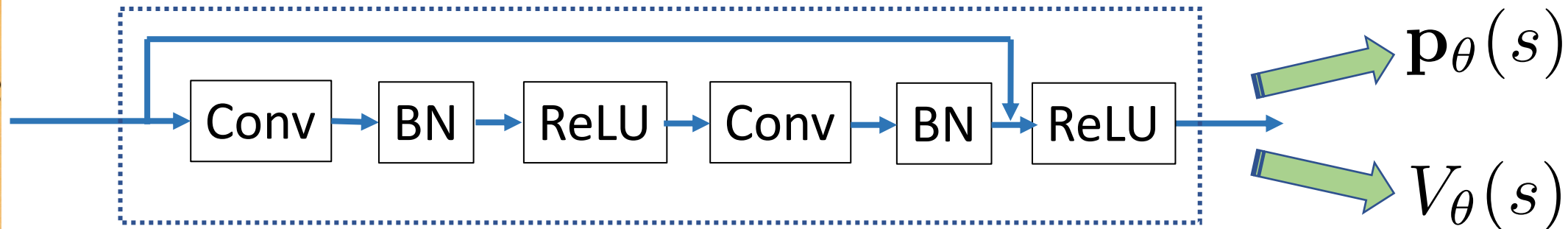
# Generate Self-play Games



# Update Models



$S$



Player situation  
at time 0

Opponent situation  
at time 0

Player situation at t=-7

Color to play

Input features (19x19x17):  $(X, Y, X_{-1}, Y_{-1}, \dots, X_{-7}, Y_{-7}, C)$

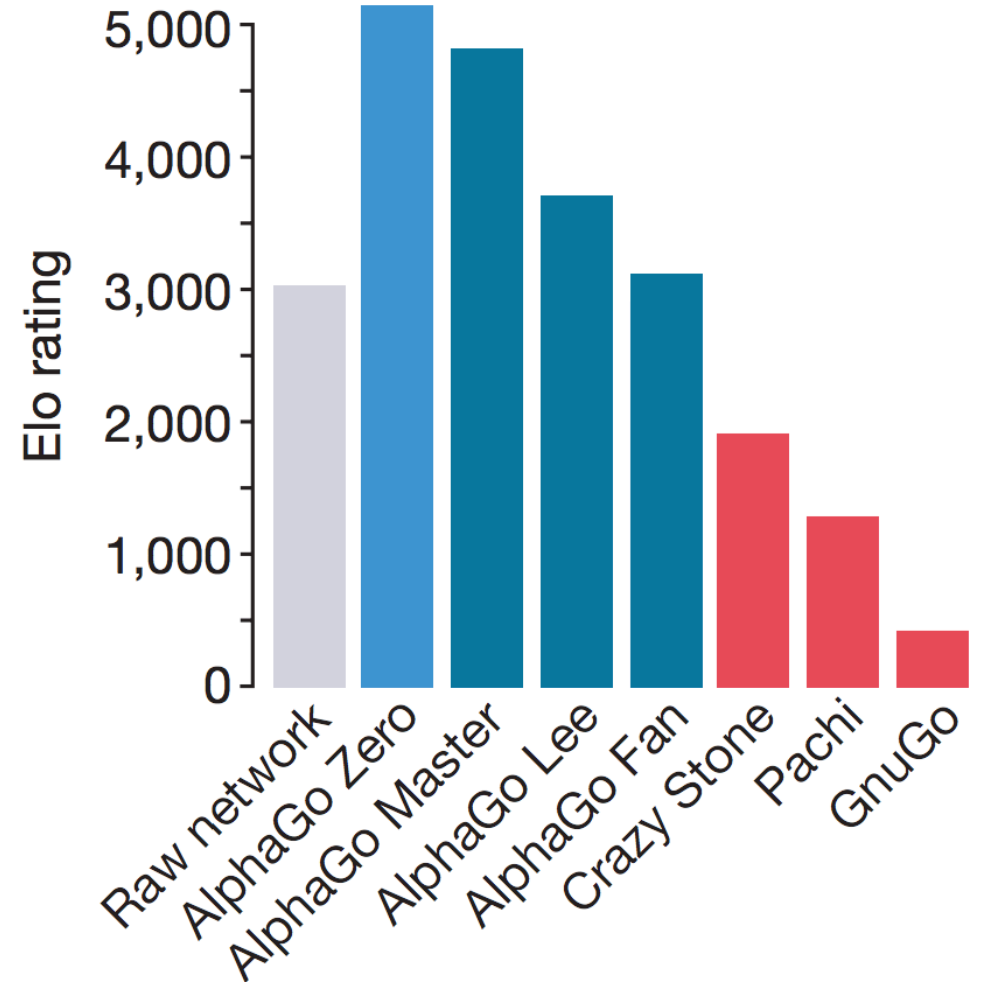
Objective:

$$J(\theta) = (z - V_{\theta})^2 - \pi^T \log \mathbf{p}_{\theta} + c \|\theta\|^2$$

$(s, \pi, z)$

# AlphaGo Zero Strength

- 3 days version
  - 4.9M Games, 1600 rollouts/move
  - 20 block ResNet
  - Defeat AlphaGo Lee.
- 40 days version
  - 29M Games, 1600 rollouts/move
  - 40 blocks ResNet.
  - Defeat AlphaGo Master by 89:11

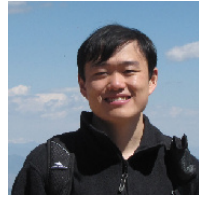


# The Mystery of AlphaZero

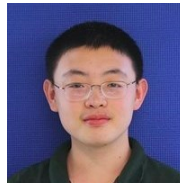
- Mystery
  - Is the proposed algorithm really universal?
  - Is the bot almighty? Is there any weakness in the trained bot?
- Lack of Ablation Studies
  - What factor is critical for the performance?
  - Is the algorithm robust to random initialization and changes of hyper parameters?
  - Any adversarial samples?

**Impressive Results, No code, No model**

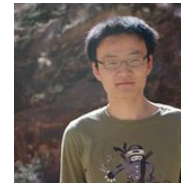
# ELF OpenGo



Yuandong Tian



Jerry Ma\*



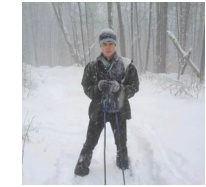
Qucheng Gong\*



Shubho Sengupta\*



Zhuoyuan Chen



James Pinkerton



Larry Zitnick

- System can be trained with 2000 GPUs in 2 weeks (20 block version)
- Superhuman performance against professional players and strong bots.
- Abundant ablation analysis.

pytorch / ELF

Unwatch 174 Unstar 2,842 Fork 472

Code Issues 36 Pull requests 3 Projects 0 Wiki Security Insights Settings Intern Dashboard

ELF: a platform for game research with AlphaGoZero/AlphaZero reimplementation

reinforcement-learning alphago-zero rl rl-environment alpha-zero go Manage topics

67 commits 11 branches 5 releases 1 environment 5 contributors View license

**We open source the code and the pre-trained model for the Go and ML community**

# ELF OpenGo Performance

## Vs top professional players

Name (rank)	ELO (world rank)	Result
Kim Ji-seok	3590 (#3)	5-0
Shin Jin-seo	3570 (#5)	5-0
Park Yeonghun	3481 (#23)	5-0
Choi Cheolhan	3466 (#30)	5-0

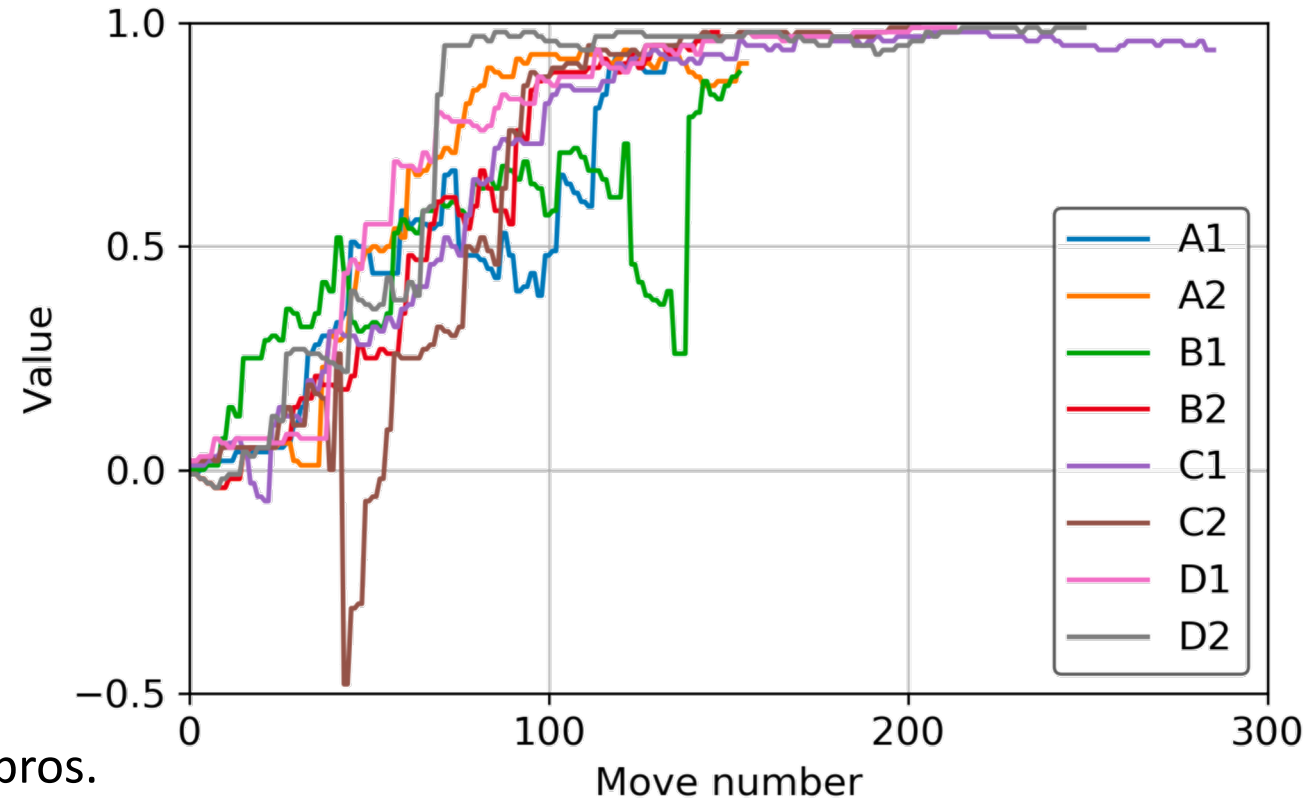
Single GPU, 80k rollouts, 50 seconds  
Offer unlimited thinking time for the players

## Vs professional players

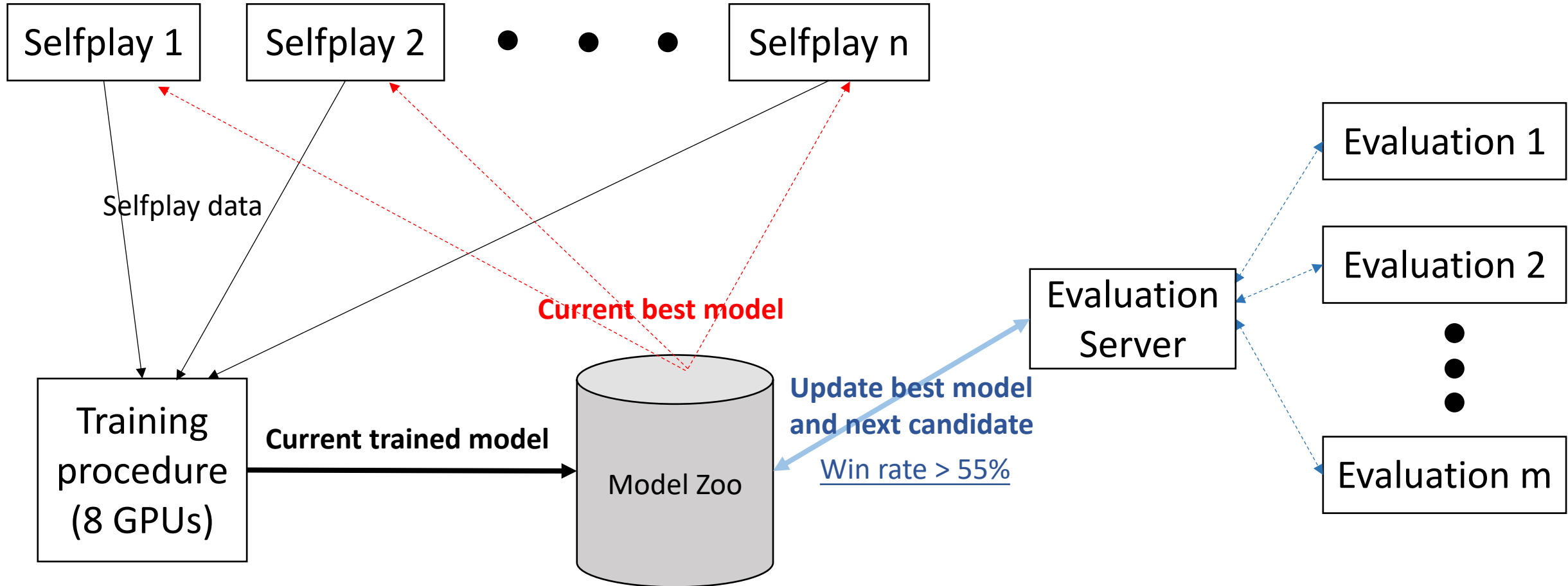
Single GPU, 2k rollouts, 27-0 against Taiwanese pros.

## Vs strong bot (LeelaZero)

[\[158603eb\]](#), 192x15, Apr. 25, 2018]: 980 wins, 18 losses (98.2%)

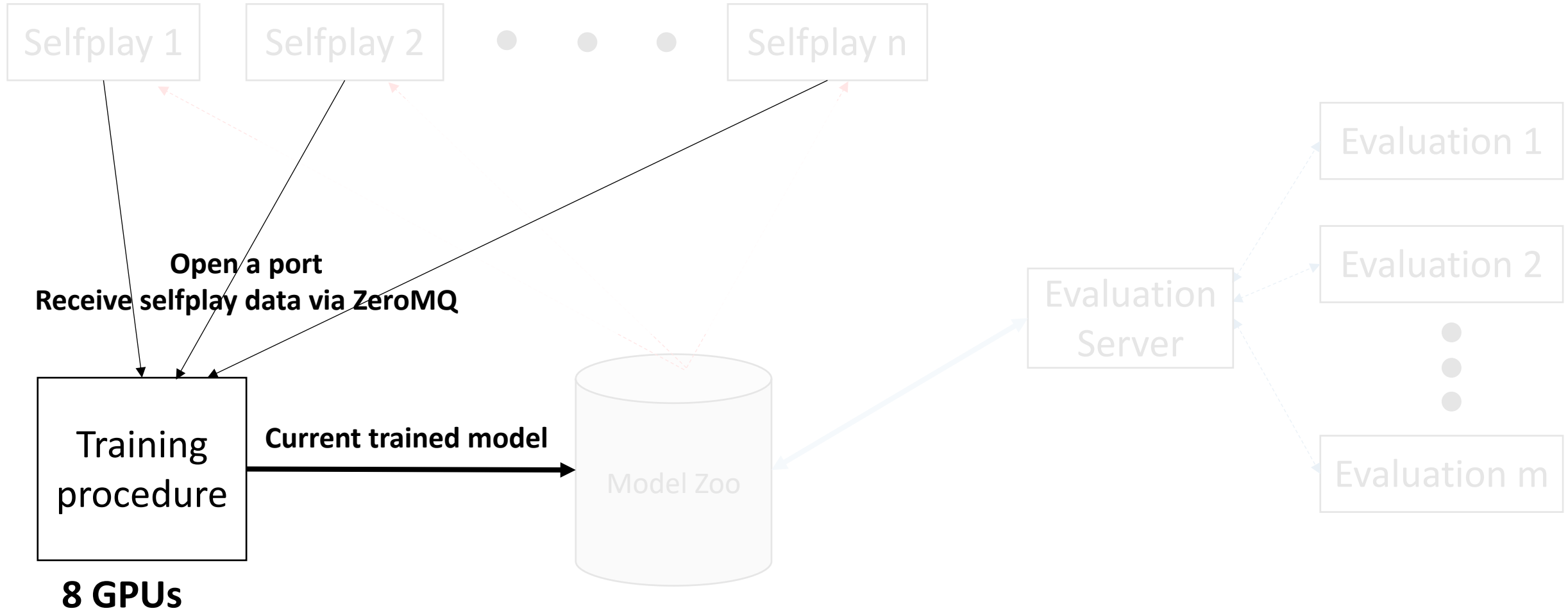


# Distributed ELF (version 1, AlphaGoZero)

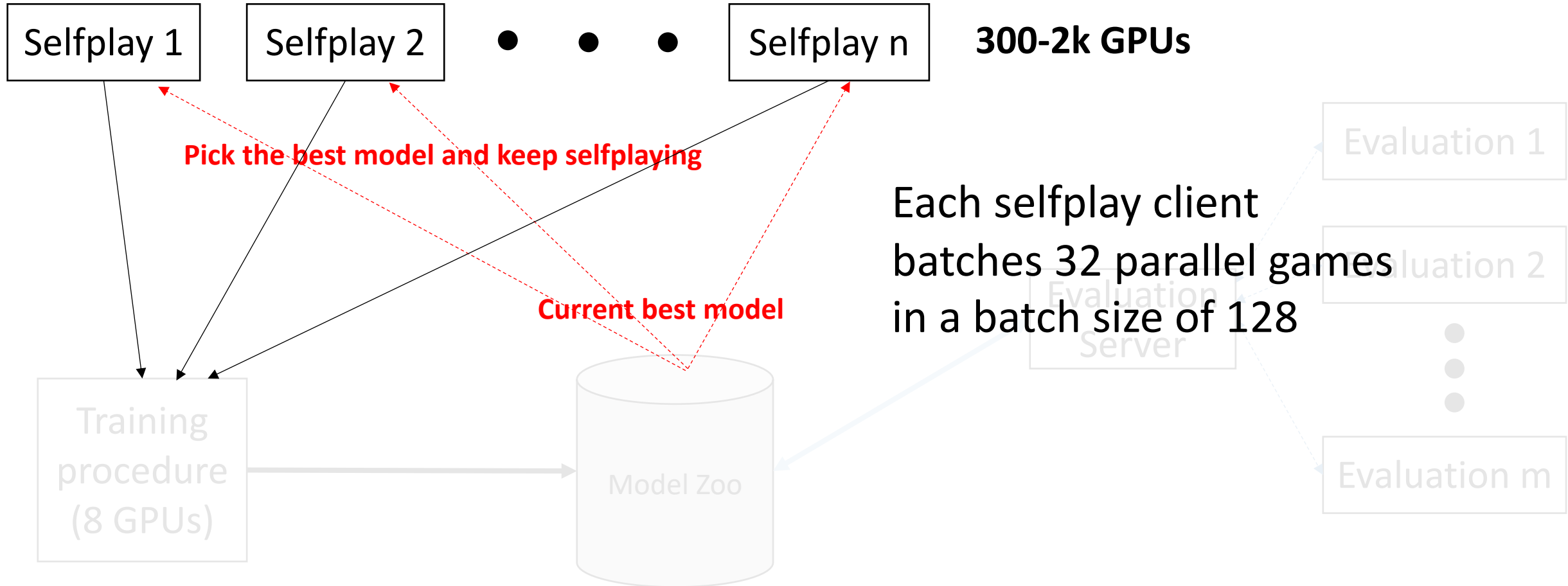




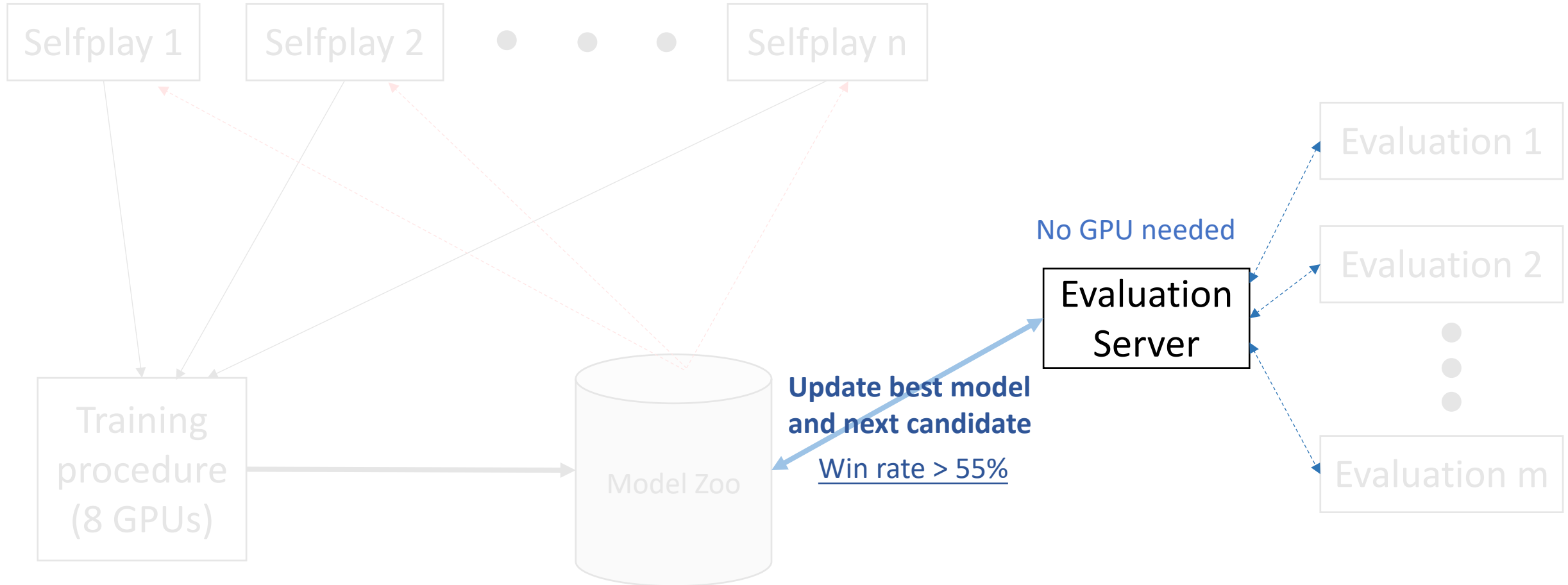
# Distributed ELF (version 1)



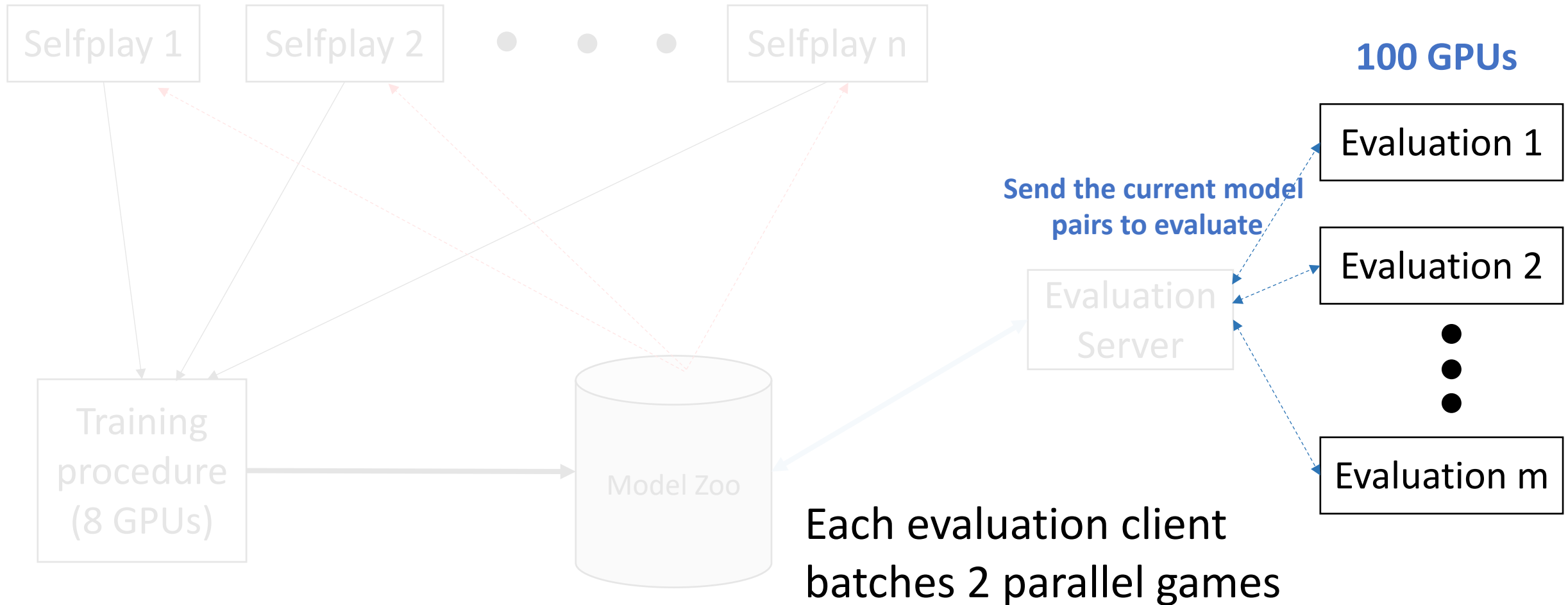
# Distributed ELF (version 1)



# Distributed ELF (version 1)



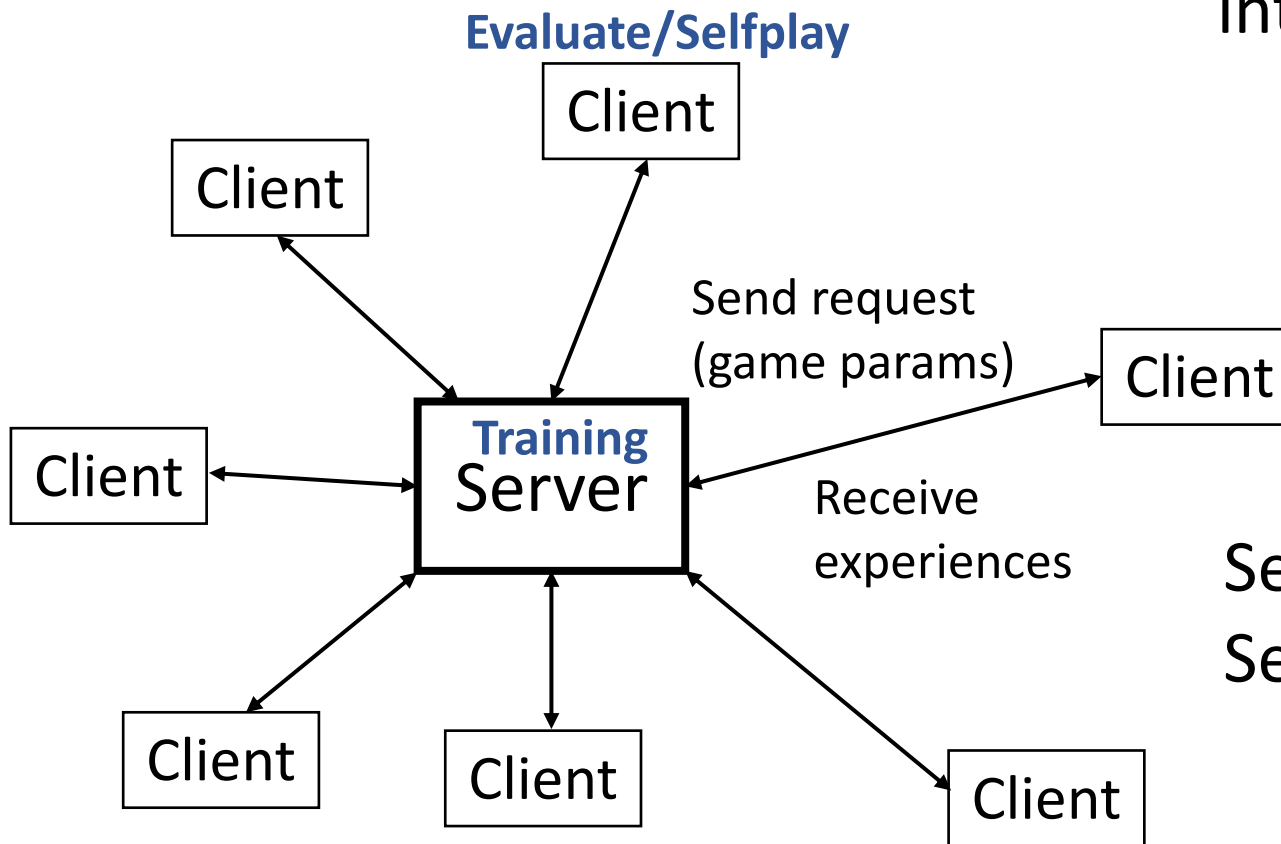
# Distributed ELF (version 1)



# Distributed ELF (v2)

Putting AlphaGoZero and AlphaZero into the same framework

AlphaGoZero (more synchronization)  
AlphaZero (less synchronization)



Server controls synchronization  
Server also does training.

# Next Step: RL Assembly

- Backbone infrastructure for ongoing projects (Hanabi, Bridge, etc)
- Reimplementation of SoTA off-policy RL methods like Ape-X and R2D2
- Incorporate OpenGo and SoTA implementation of MCTS.
- Efficient on single machine (SoTA training FPS so far)

**Open source soon**

# Frame Per Second (FPS) on Atari Games

ReLA: 12.5 KFPS

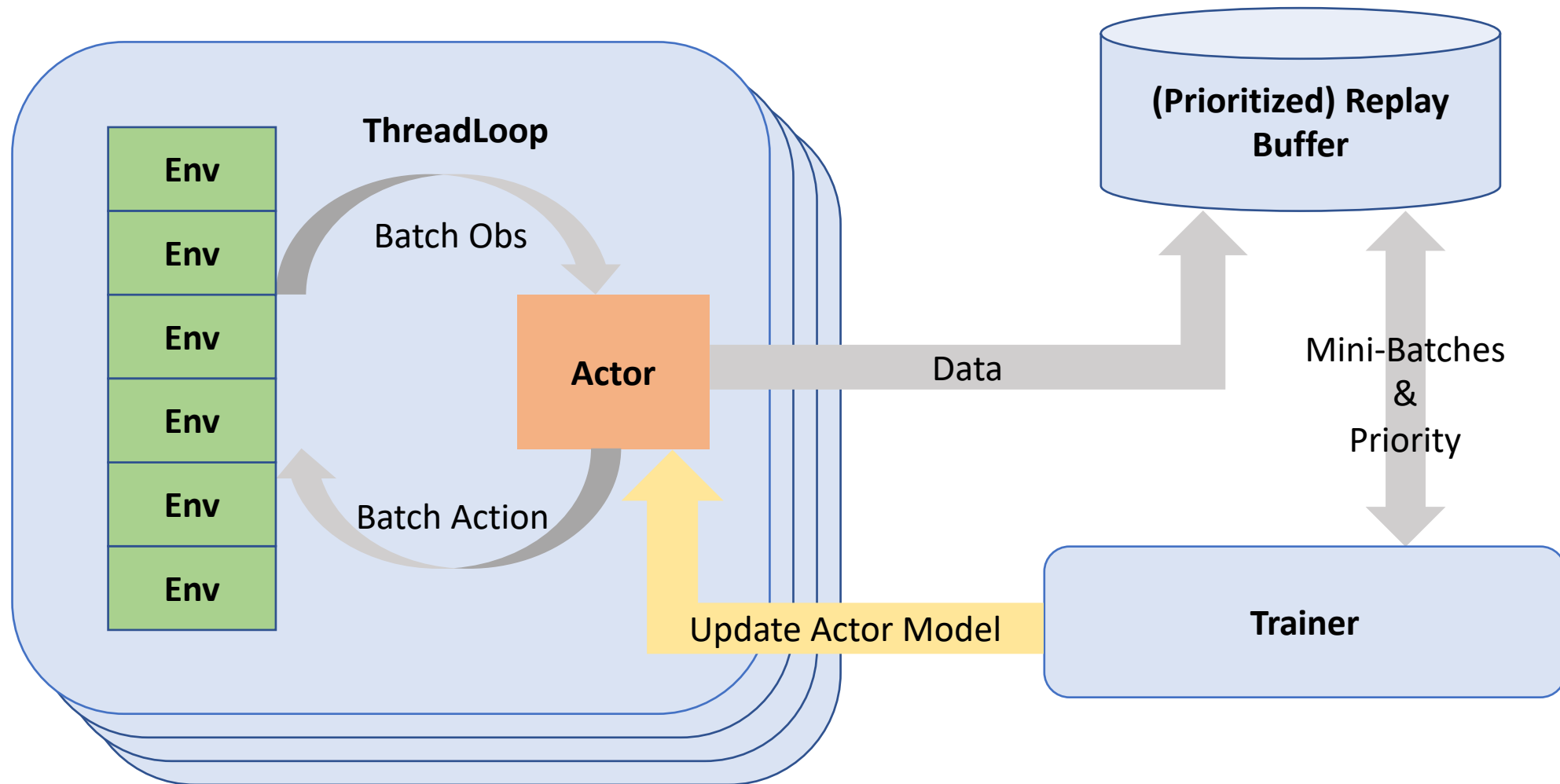
*using 40 CPU cores + 2 GPU (P100) on a single machine*

Ape-X: 12.5 KFPS

*using 360 CPU cores + 1 GPU (distributed system)*

- ReLA is GPU bound. Performance is better with more GPUs
- A few more improvements to achieve better performance when releasing.

# Architecture





# User Interface (API)

```
1  env = reLa.VectorEnv()
2  for _ in range(num_env_per_thread):
3      game = create_atari(...)
4      env.append(game)
5  actor = reLa.DQNActor(...)
6  thread = reLa.ThreadLoop(actor, env)
```

All objects (env, agent, replay buffer, etc) are created & configured in Python

# User Interface (API)

```
class ApexAgent(torch.jit.ScriptModule):  
    @torch.jit.script_method  
    def td_err(self, obs: Dict[str, torch.Tensor], ...) -> torch.Tensor:  
        online_q = self.online_net(obs)  
        online_qa = online_q.gather(1, action.unsqueeze(1)).squeeze(1)  
  
        next_obs = self.greedy_act(next_obs)  
        bootstrap_q = self.target_net(next_obs)  
        bootstrap_qa = bootstrap_q.gather(1, next_a.unsqueeze(1)).squeeze(1)  
        target = reward + bootstrap * (self.gamma ** self.multi_step) * bootstrap_qa  
        return target.detach() - online_qa
```

Model is written in **Python** with **PyTorch's TorchScript**,  
and executed in **C++** with multi-threading for maximum throughput.

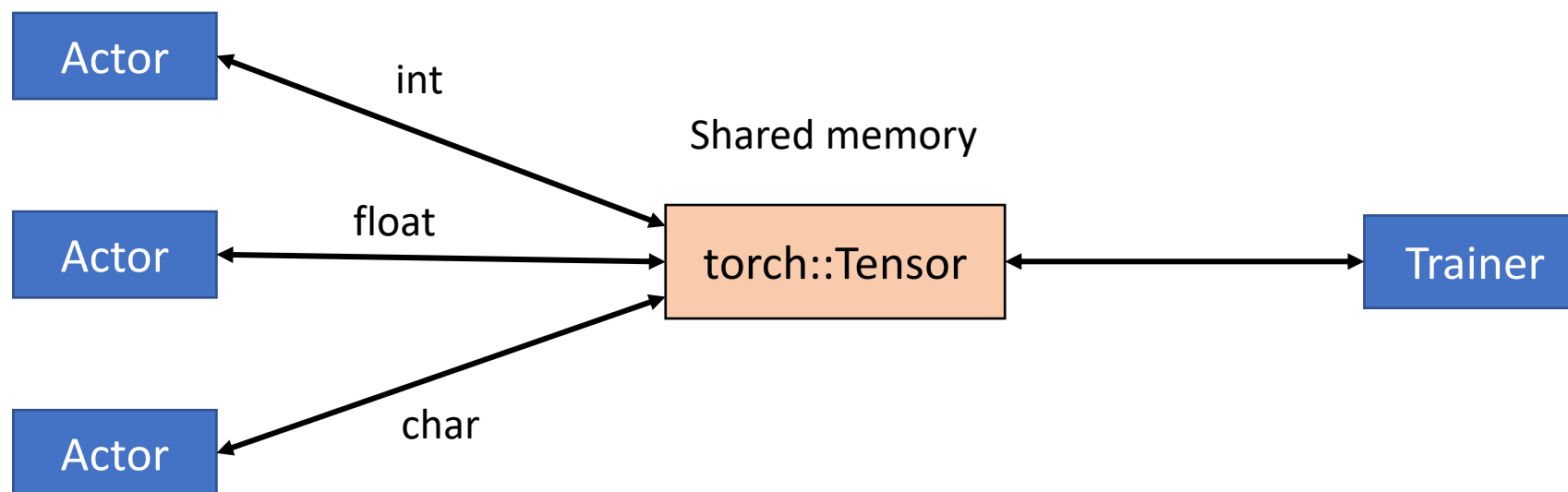
# Native integration with PyTorch C++ API

- Simple/Intuitive manipulations of PyTorch tensors in C++
  - Same as/Similar to Python Interface
  - No extra library needed for operations like downsample/upsampling.

```
1 torch::Tensor s = getObservation();
2 s = s.view({1, 3, height, width});
3 // rescale the image
4 s = torch::upsample_bilinear2d(s, {sHeight, sWidth}, true);
5 s = s.view({3, sHeight, sWidth});
6 // convert to grey scale
7 s = 0.21 * s[0] + 0.72 * s[1] + 0.07 * s[2];
```

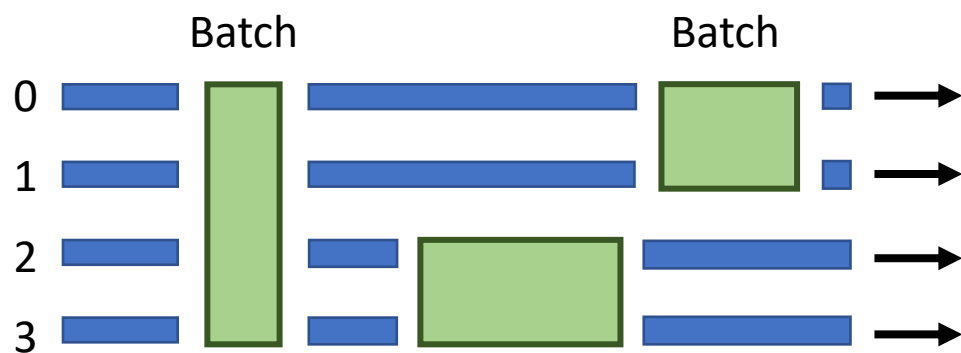
# Native integration with PyTorch C++ API

- Easier communication between threads/processes via Tensor.
  - No extra copy when sending data from/to environments.

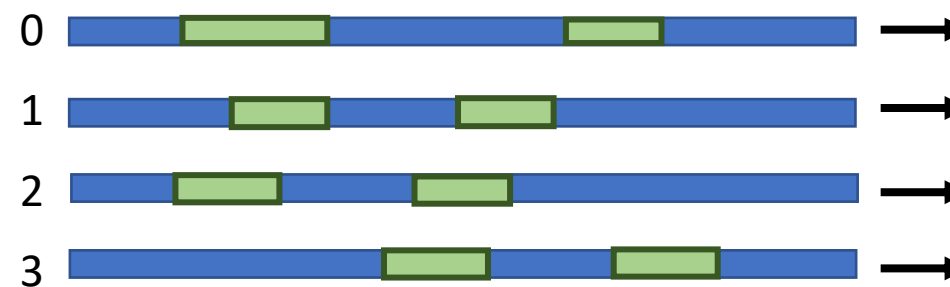


# Native integration with PyTorch C++ API

- Simultaneous network forwarding at different threads
  - Python GIL becomes irrelevant.
  - No need to block the environment
    - good for simple environments like Go, Bridge, Hanabi and others.



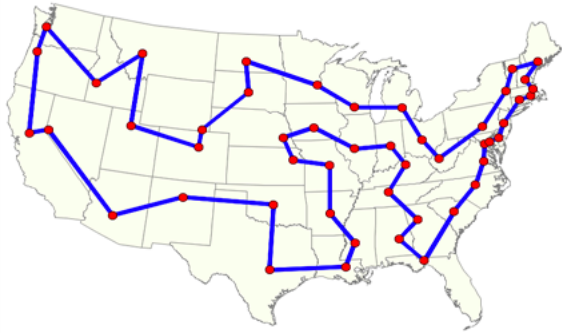
Game threads are blocked  
Waiting until python side to reply



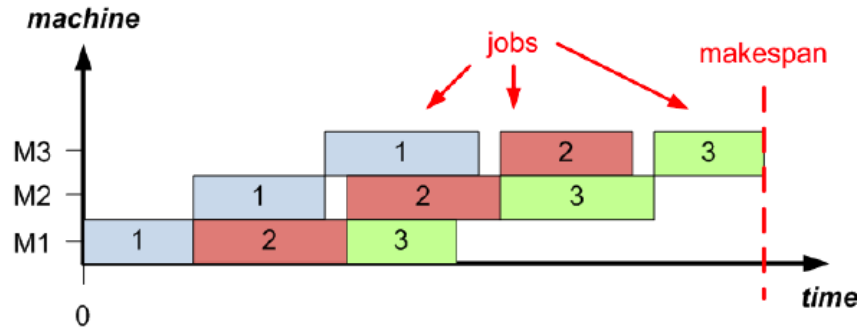
Game threads calling PyTorch  
C++ API directly

# Learning Hand-tuned Heuristics with RL/ML

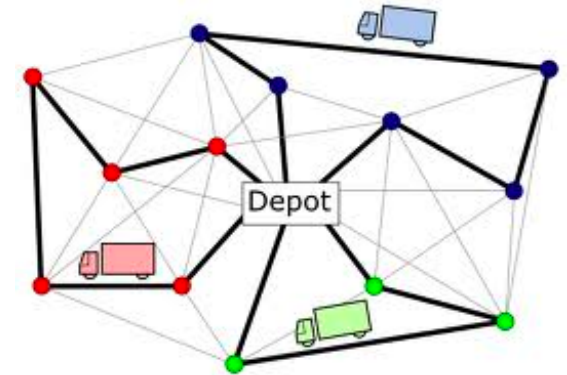
# Combinatorial optimization



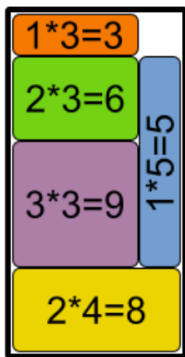
Travel Salesman Problem



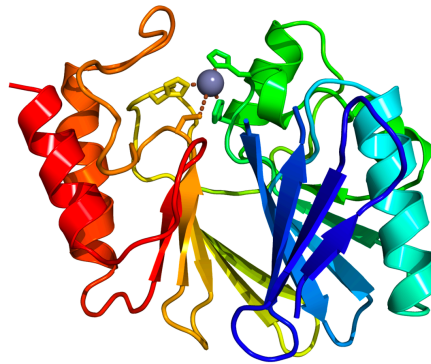
Job Scheduling



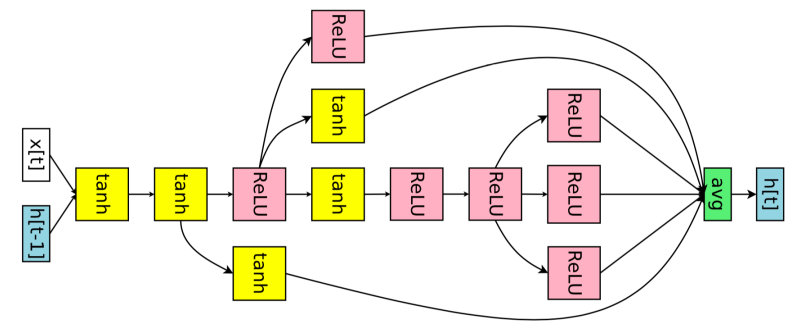
Vehicle Routing



Bin Packing



Protein Folding



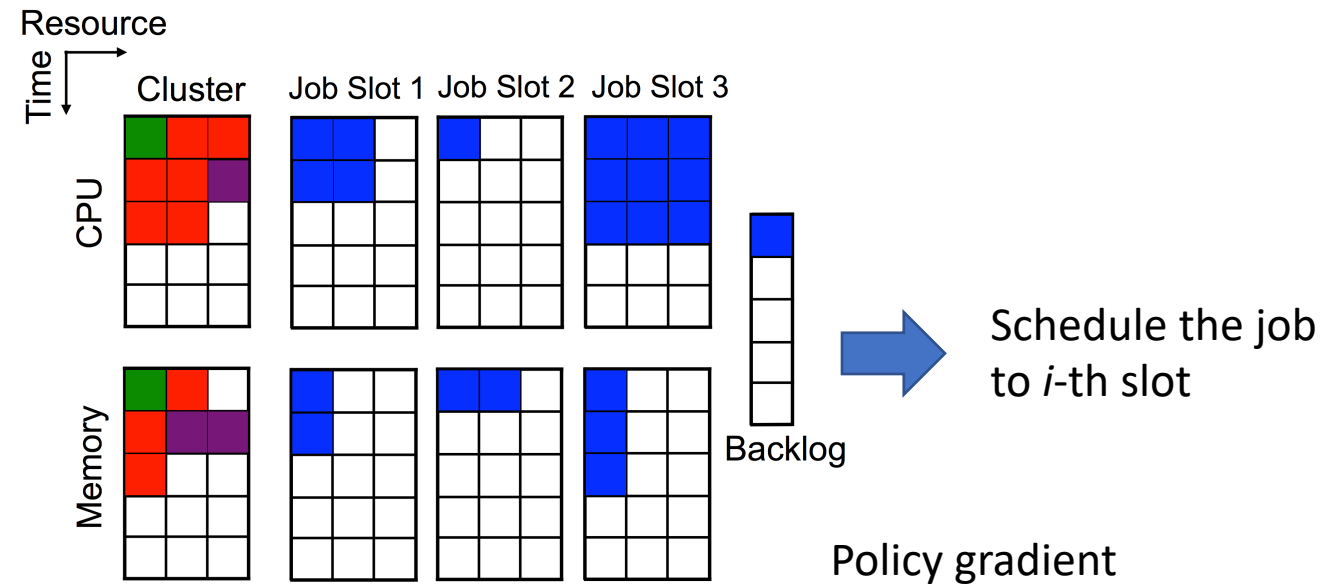
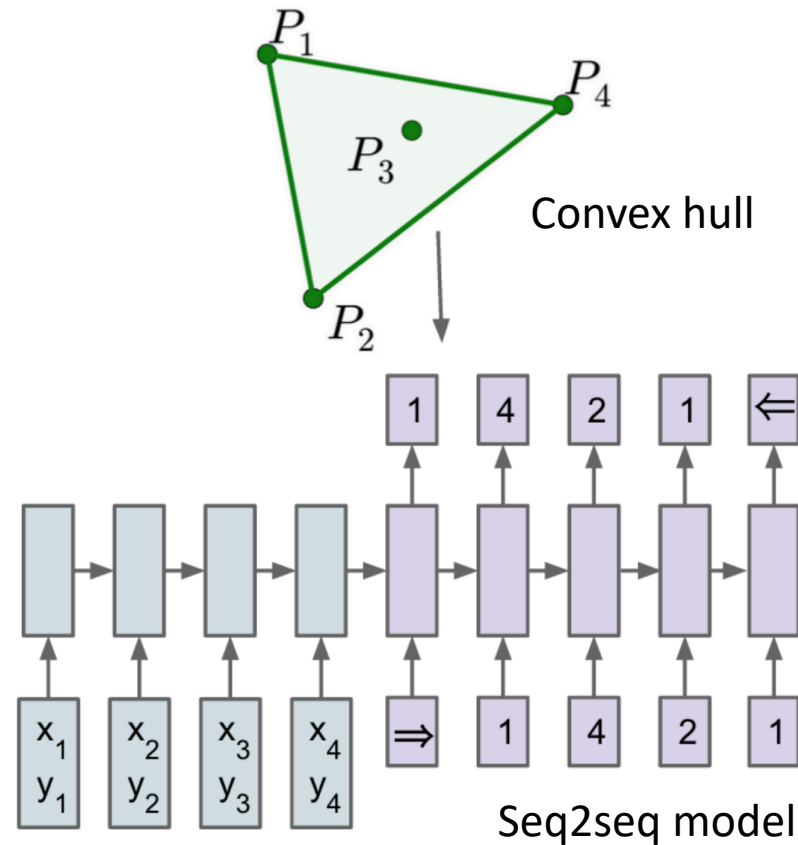
Model-Search

# Wait...What?

- These are NP-hard problems.
  - No good algorithm unless  $P = NP$
- These guarantees are worst-case ones.
  - To prove a lower-bound, construct an adversarial example to fail the algorithm
- For specific distribution, there might be better heuristics.
  - Human heuristics are good but may not be suitable for everything

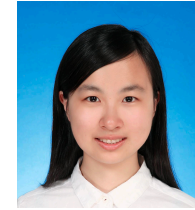


# Direct predicting combinatorial solutions

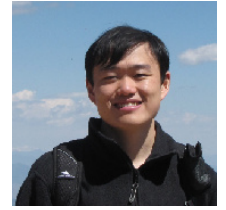


[H. Mao et al, Resource Management with Deep Reinforcement Learning, ACM Workshop on Hot Topics in Networks, 2016]

# Local Rewriting Framework



Xinyun Chen

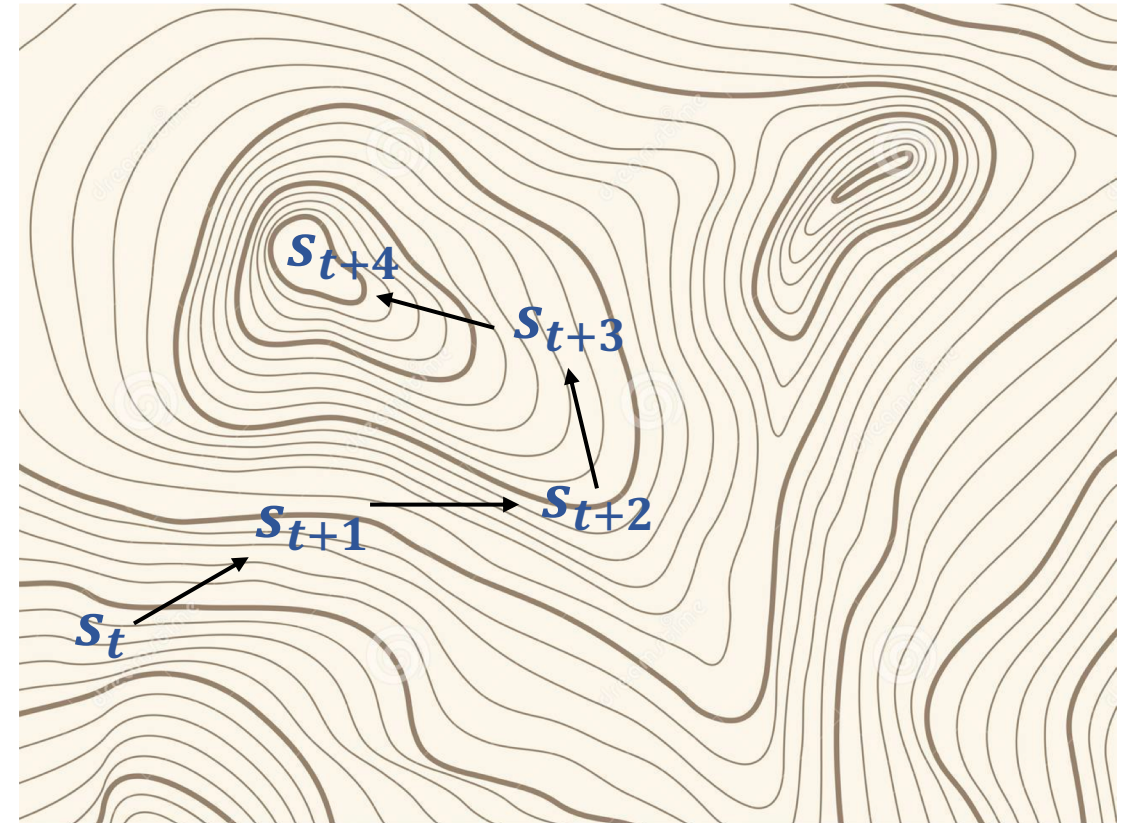


Yuandong Tian

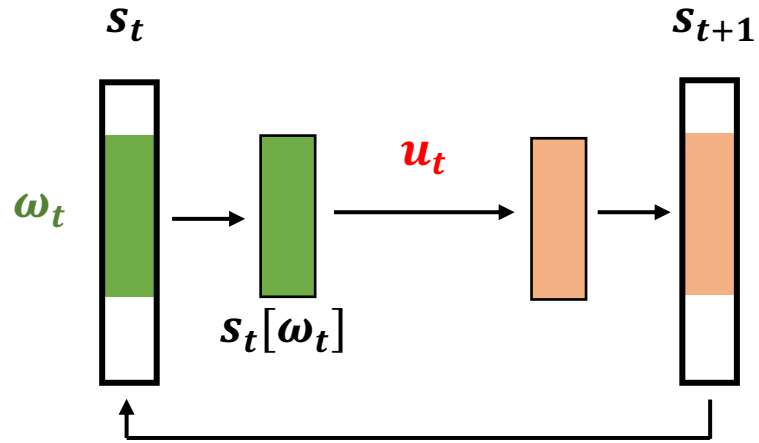
Code: <https://github.com/facebookresearch/neural-rewriter>

A learned “gradient descent” that  
starts from a feasible solution  
iteratively converges to a good solution

How to learn it?



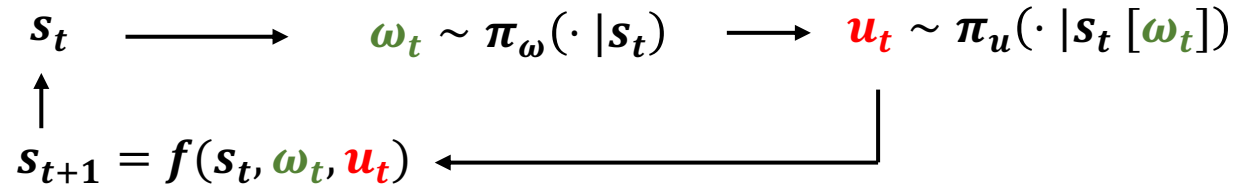
# Local Rewriting Framework



Current State  
(i.e. Solution)

Region-Picker

Rule-Picker



# Q-Actor-Critic Training

How to train two policies  $\pi_{\omega}(\cdot | \mathbf{s}_t)$  and  $\pi_u(\cdot | \mathbf{s}_t [\boldsymbol{\omega}_t])$ ?

Learn Q to fit cumulative rewards:

$$L_{\omega}(\theta) = \frac{1}{T} \sum_{t=0}^{T-1} \left( \sum_{t'=t}^{T-1} \gamma^{t'-t} r(s'_{t'}, (\omega'_{t'}, u'_{t'})) - Q(s_t, \omega_t; \theta) \right)^2$$

$\pi_{\omega}(\cdot | \mathbf{s}_t)$ : Q-learning with soft policy:

$$\pi_{\omega}(\omega_t | s_t; \theta) = \frac{\exp(Q(s_t, \omega_t; \theta))}{\sum_{\omega_t} \exp(Q(s_t, \omega_t; \theta))}$$

$\pi_u(\cdot | \mathbf{s}_t [\boldsymbol{\omega}_t])$ : Actor-Critic with learned Q:

$$L_u(\phi) = - \sum_{t=0}^{T-1} \Delta(s_t, (\omega_t, u_t)) \log \pi_u(u_t | s_t[\omega_t]; \phi)$$

Advantage:

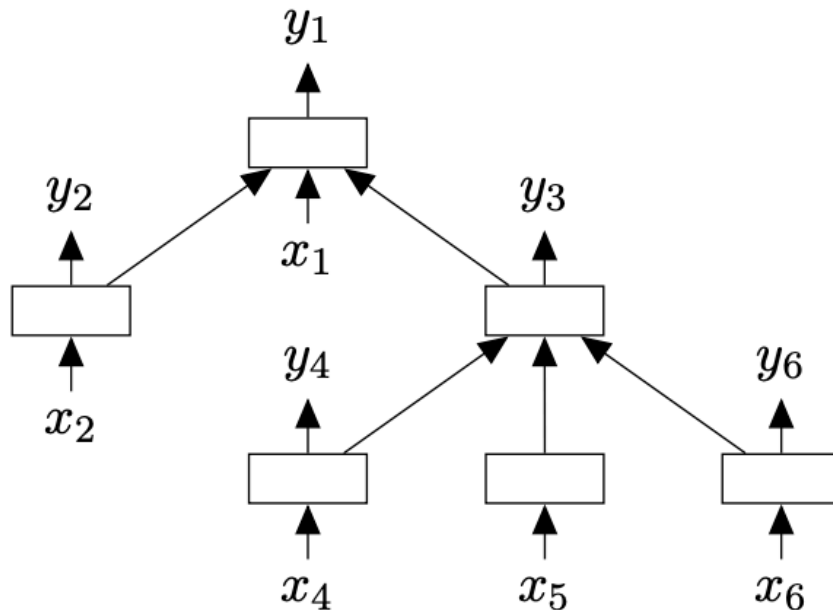
$$\Delta(s_t, (\omega_t, u_t)) \equiv \sum_{t'=t}^{T-1} \gamma^{t'-t} r(s'_{t'}, (\omega'_{t'}, u'_{t'})) - Q(s_t, \omega_t; \theta)$$

# How to encode Structure Data

## Child-Sum LSTM

$$y_1 = f(y_2, y_3, x_1)$$

$f$  can be very complicated:



$$\tilde{h}_j = \sum_{k \in C(j)} h_k,$$

$$i_j = \sigma \left( W^{(i)} x_j + U^{(i)} \tilde{h}_j + b^{(i)} \right),$$

$$f_{jk} = \sigma \left( W^{(f)} x_j + U^{(f)} h_k + b^{(f)} \right),$$

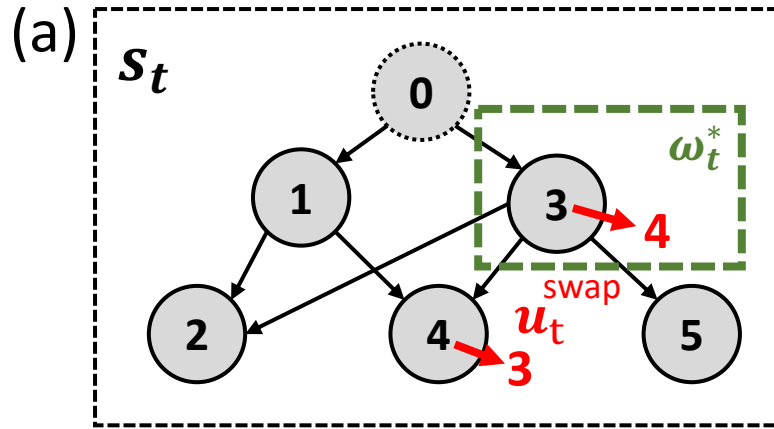
$$o_j = \sigma \left( W^{(o)} x_j + U^{(o)} \tilde{h}_j + b^{(o)} \right),$$

$$u_j = \tanh \left( W^{(u)} x_j + U^{(u)} \tilde{h}_j + b^{(u)} \right),$$

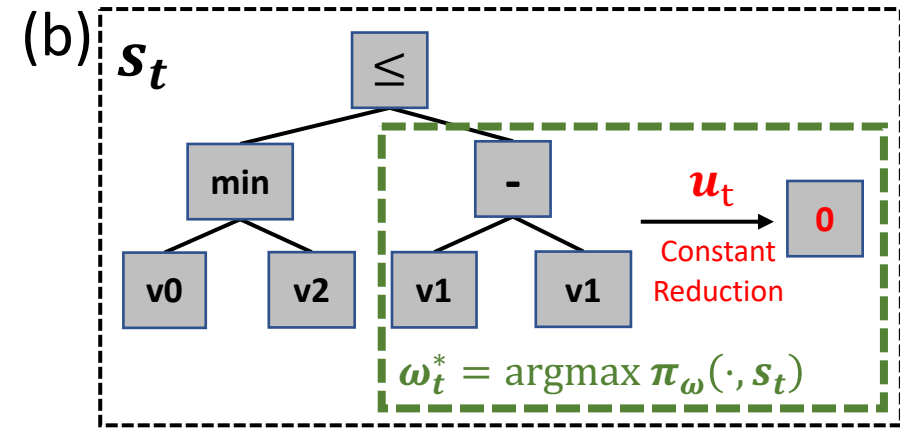
$$c_j = i_j \odot u_j + \sum_{k \in C(j)} f_{jk} \odot c_k,$$

$$h_j = o_j \odot \tanh(c_j),$$

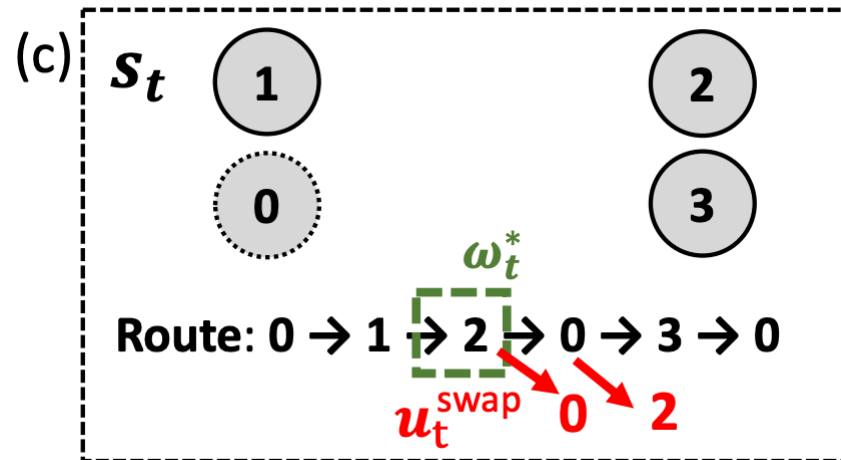
# Applications



Online Job Scheduling

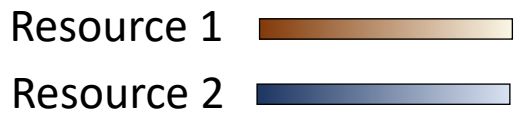
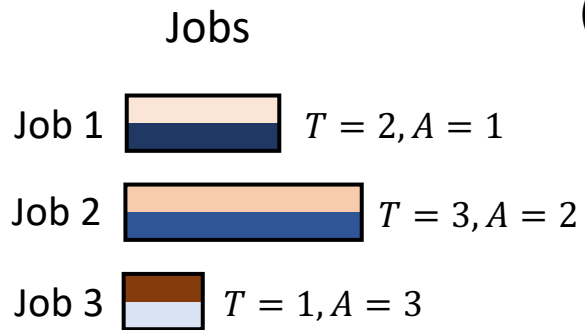


Expression Simplification

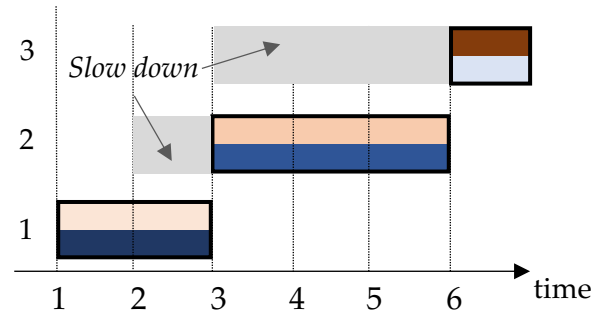


Vehicle Routing

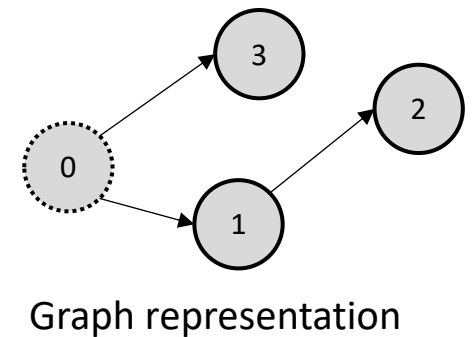
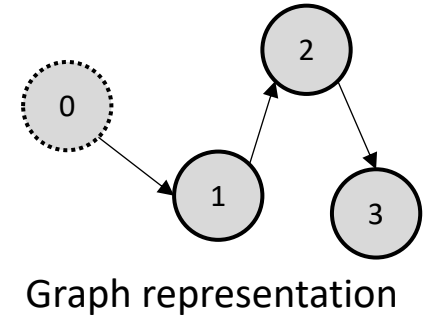
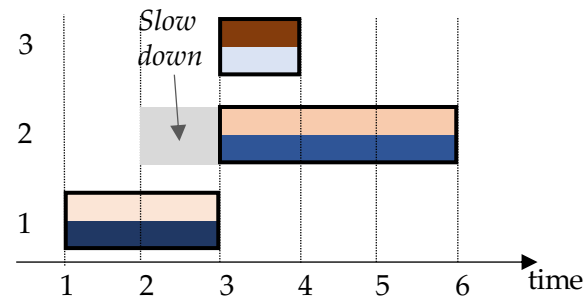
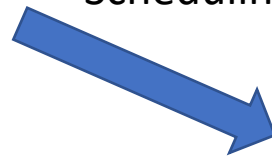
# Online Job Scheduling



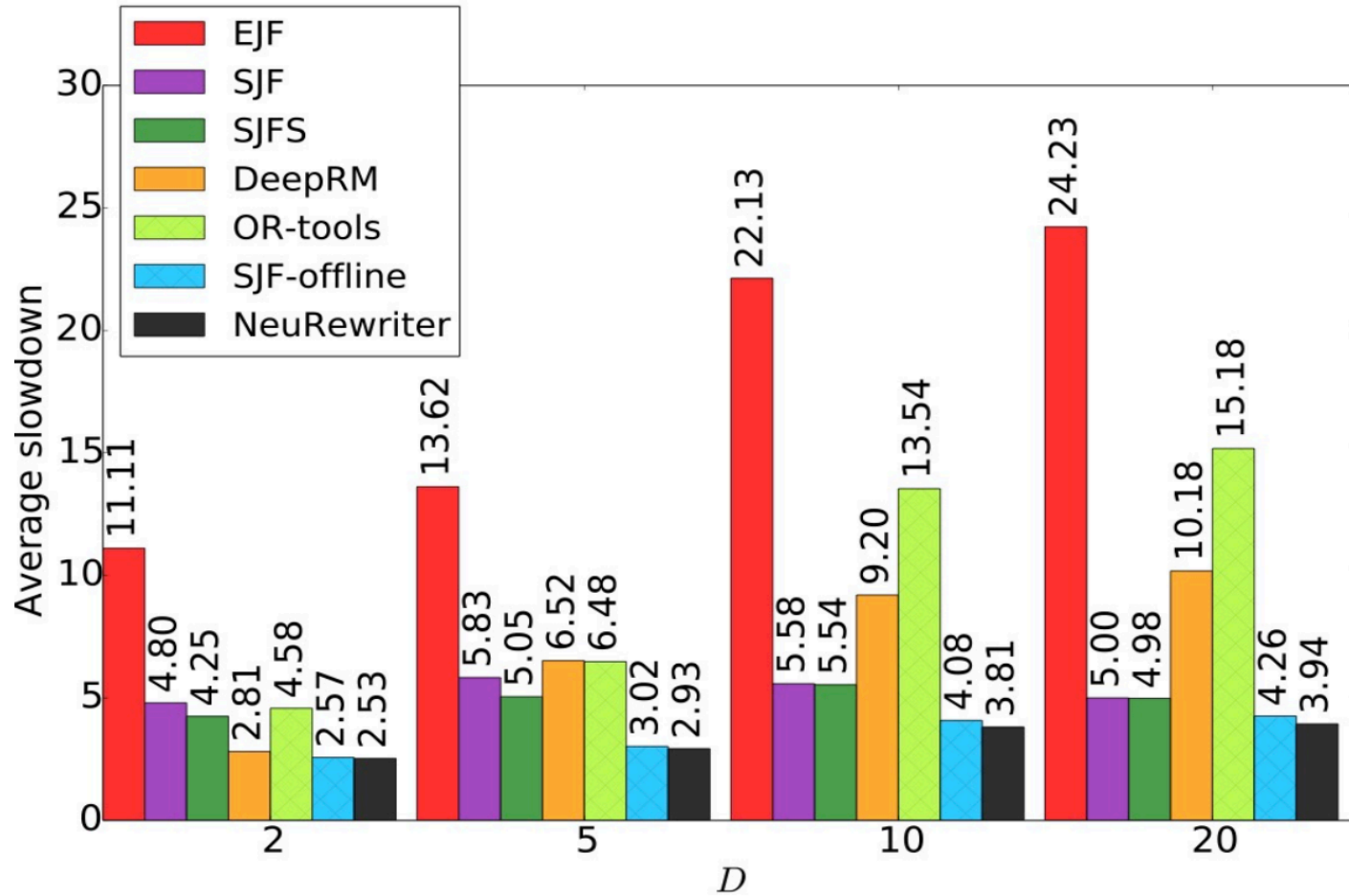
Scheduling 1  
(Sequential)



Scheduling 2



# Online Job Scheduling



Baselines:

- Earliest Job First (EJF)
- Shortest Job First (SJF)
- Shortest First Search (SJFS)
- DeepRM

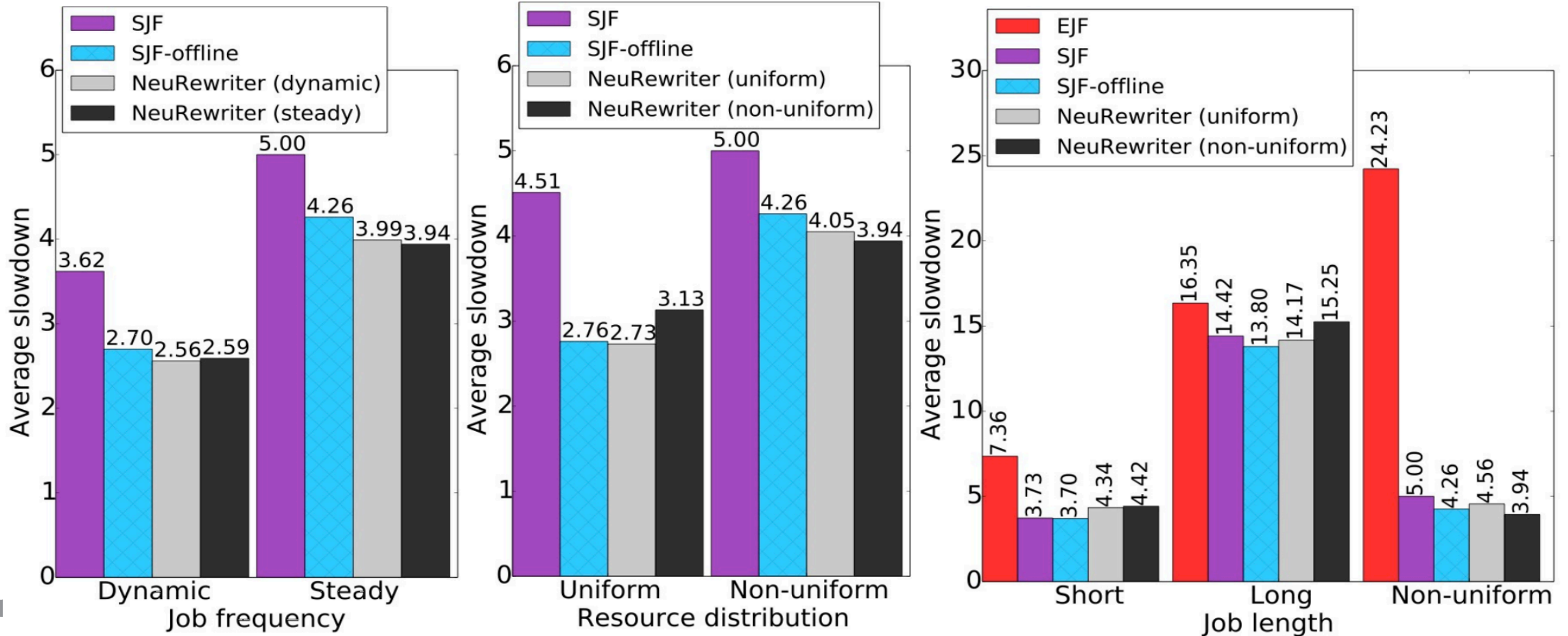
Offline baselines:

- Google OR-tools (OR-tools)
- SJF-offline

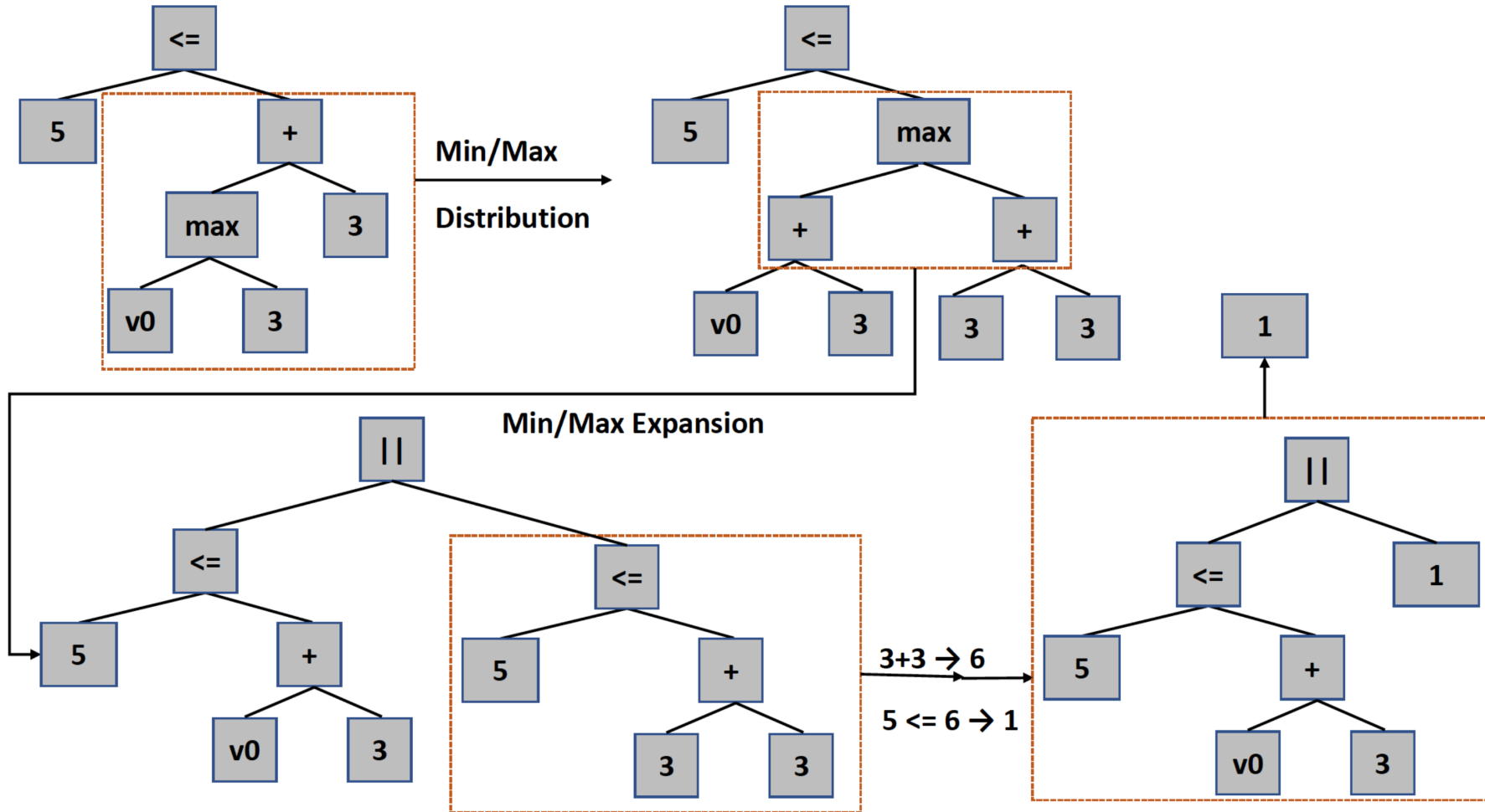


# Online Job Scheduling: Ablation Study

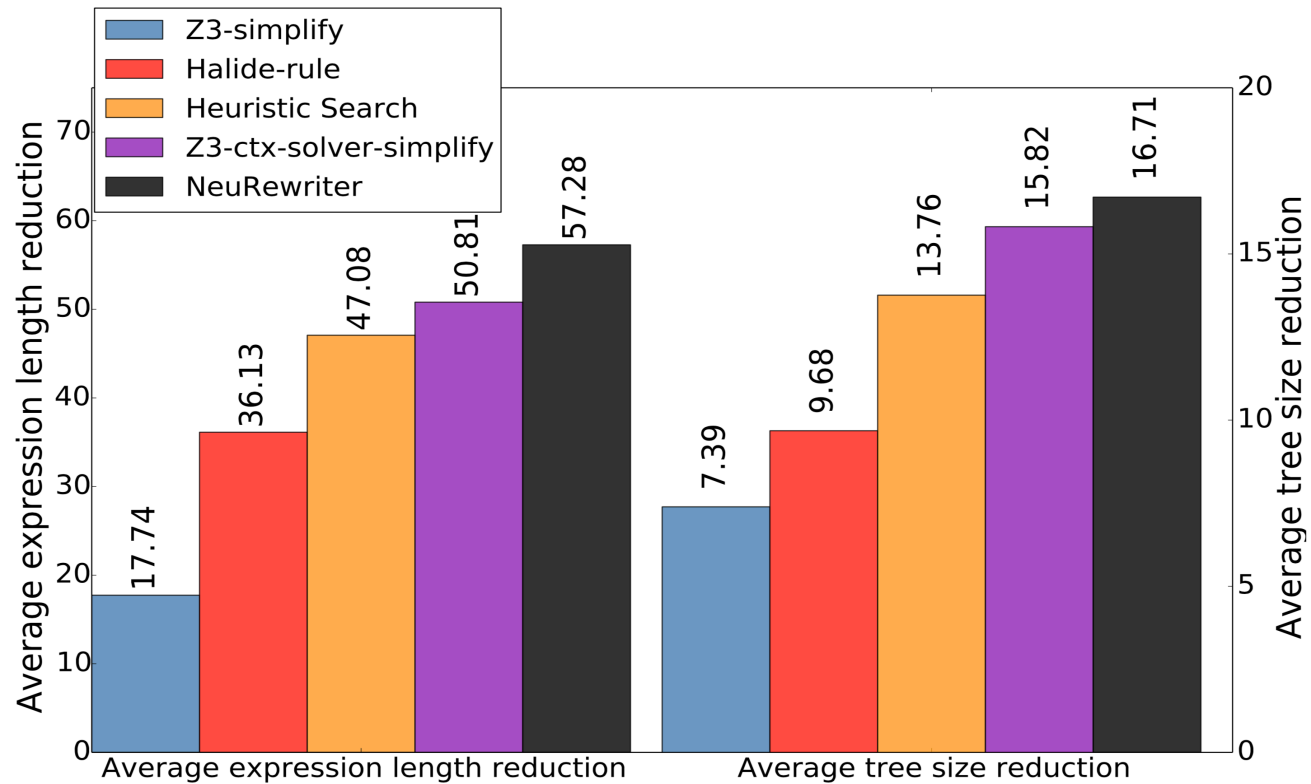
The learned model can generalize to different job distributions.



# Expression Simplification



# Expression Simplification



Baselines:

Z3-simplify

Z3-ctx-solver-simplify

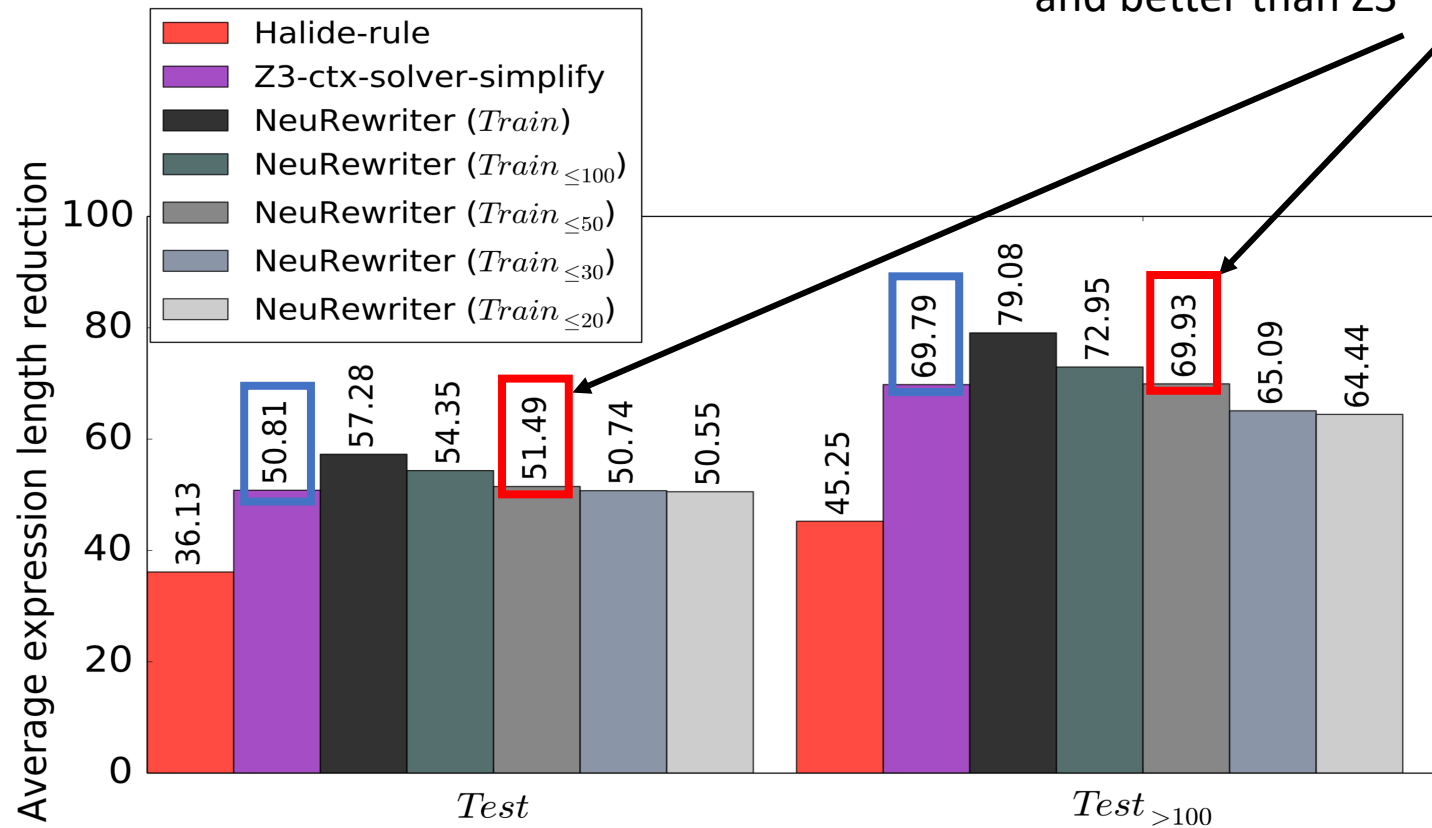
Heuristic Search

Halide rules

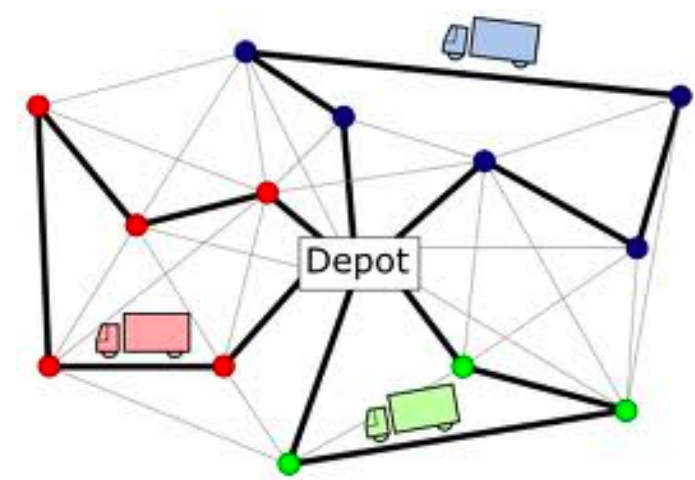
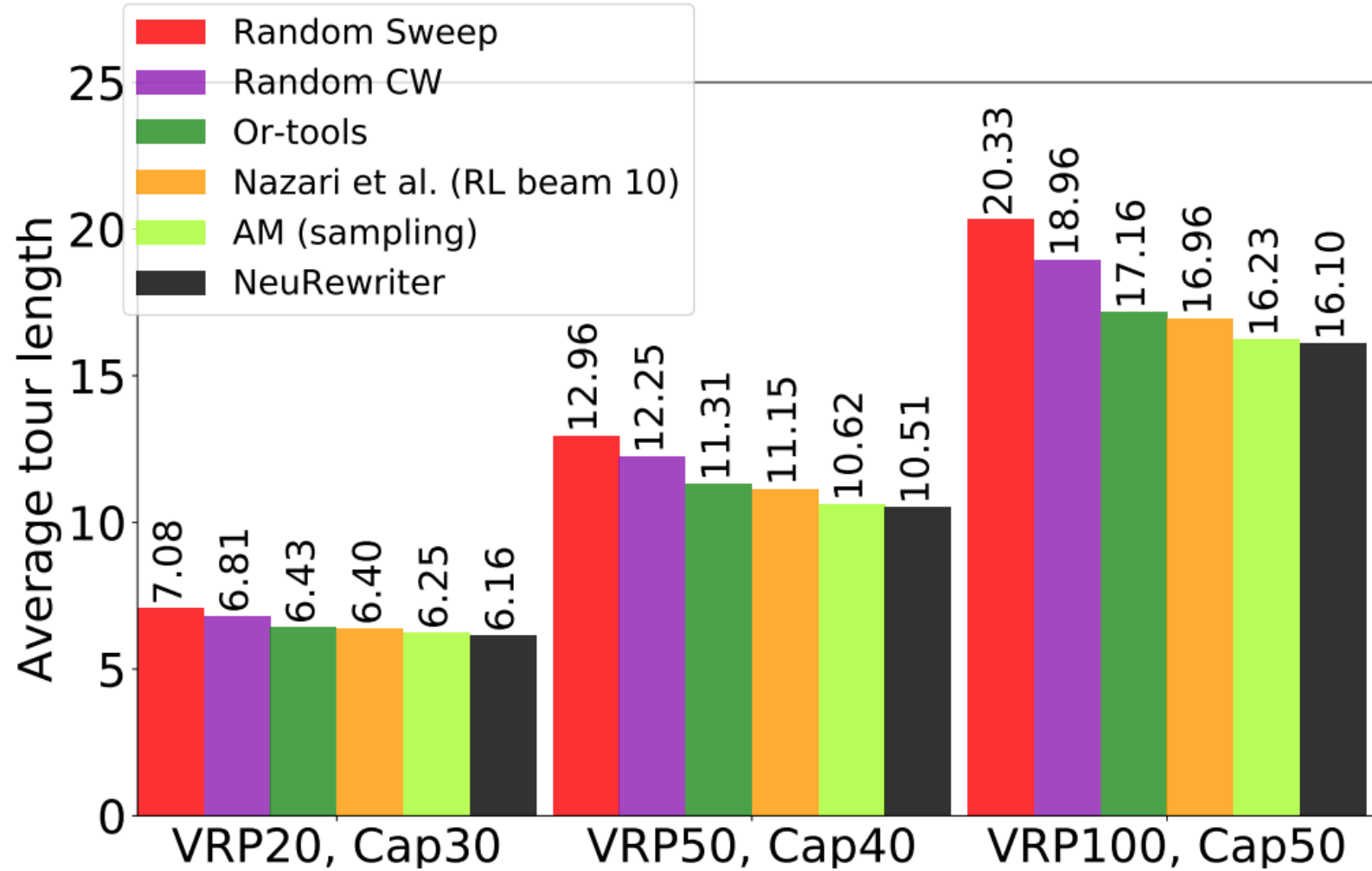
# Expression Simplification

Transfer learning still works well.

A model trained with expression length  $\leq 50$  has good performance on test set with expression length  $\geq 100$ , and better than Z3



# Capacitated Vehicle Routing



# Coda: An End-to-End Neural Program Decompiler

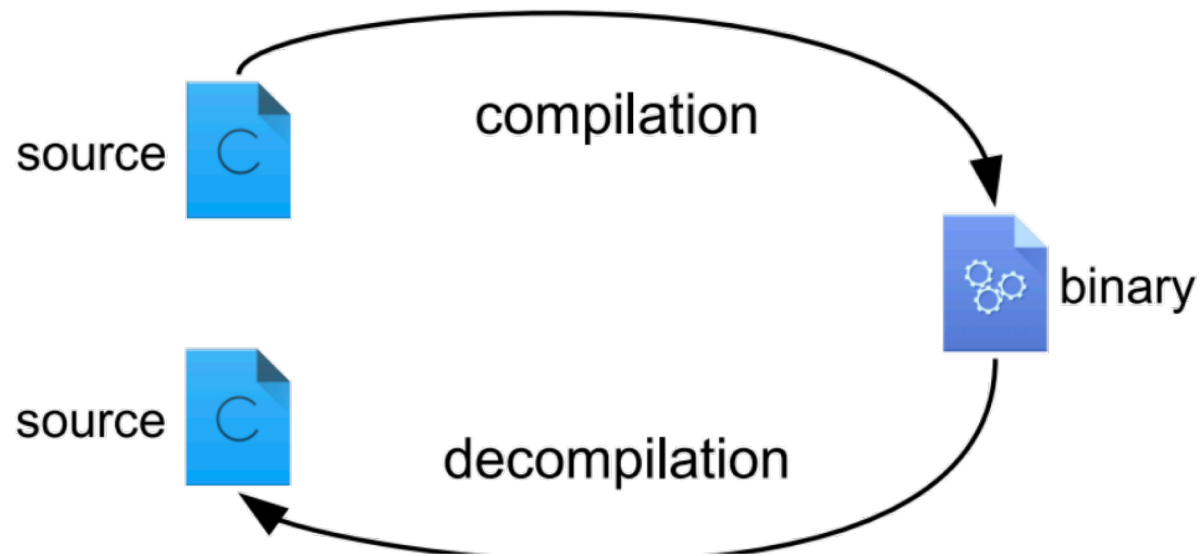
Cheng Fu<sup>1</sup>, Huili Chen<sup>1</sup>, Haolan Liu<sup>1</sup>, Xinyun Chen<sup>3</sup>, Yuandong Tian<sup>2</sup>, Farinaz Koushanfar<sup>1</sup>, Jishen Zhao<sup>1</sup>

*<sup>1</sup>UC San Diego, <sup>2</sup>Facebook AI Research, <sup>3</sup>UC Berkeley*

NeurIPS 2019

# Background: Decompilation

- Goal of Decompilation
  - From Binary Execution to High-level program language



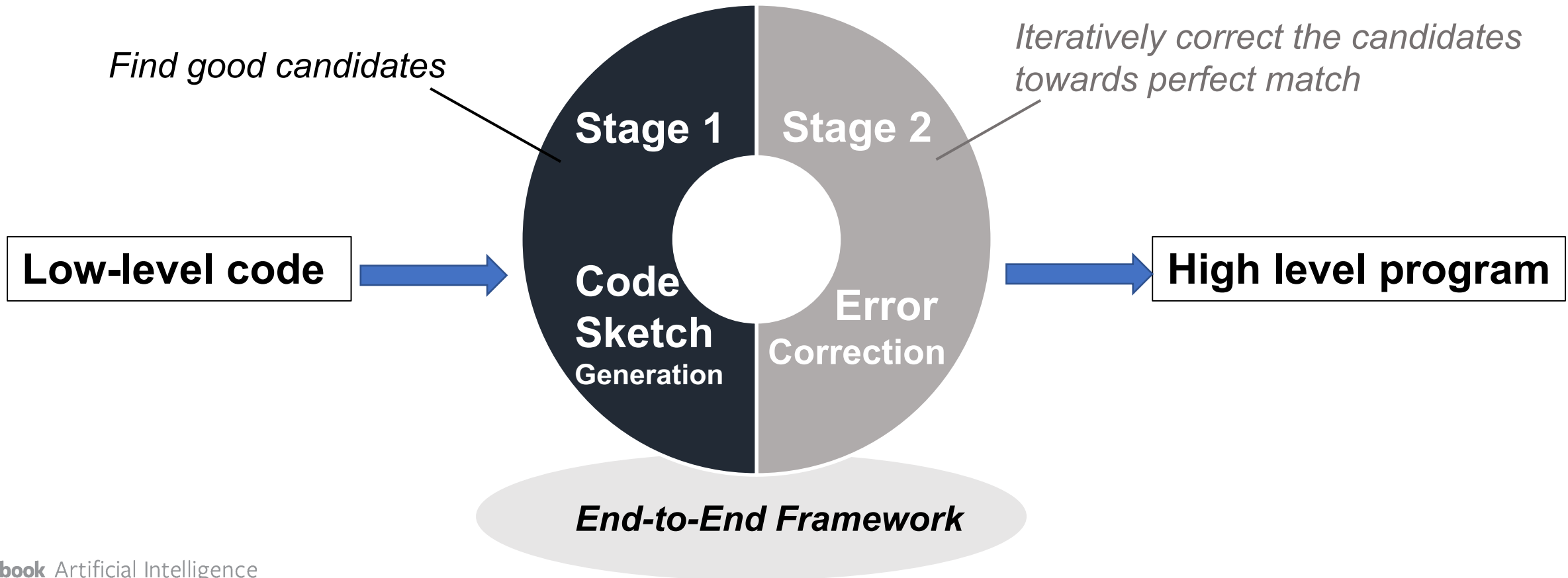
# Challenges

- Many hardware architectures (ISA): x86, MIPS, ARM
- Many Programming Languages (PL)
  - Extra Human effort to extend to the new version of the hardware architectures or programming languages
- Our goals:
  - Maintain both the functionality and semantics of the binary executables
  - Make the design process end-to-end (generalizable to various ISAs and PLs)



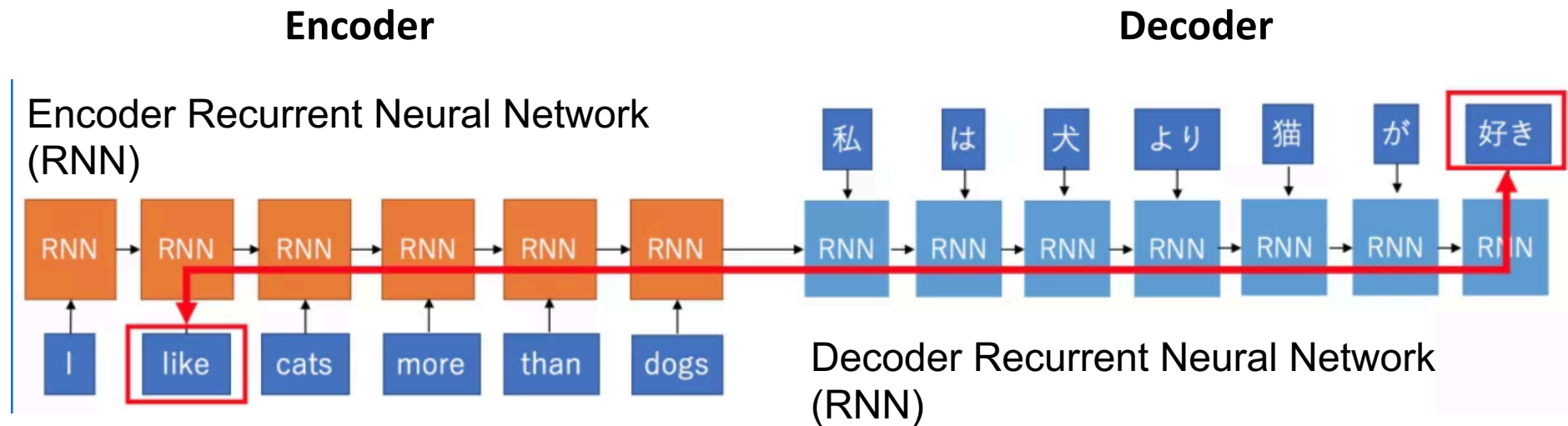
# Coda Design

*Leverage both syntax and dynamic information*

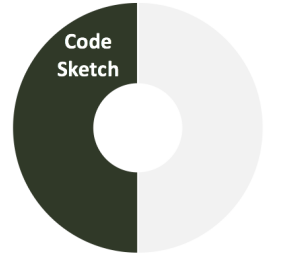


# Stage 1: Coda Sketch Generation

- Is Decompilation simply a translation problem?



**More than a translation problem!**



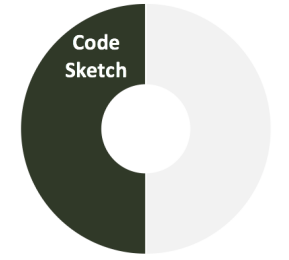
# Stage 1: Coda Sketch Generation

## • Encoder

- N-ary Tree Encoder to capture **inter** and **intra** dependencies of the low-level code.
- Opcode and its operands are encoded together
- Different encoder is used for different instruction types
  - memory (mem)
  - branch (br)
  - arithmetic (art).

# source C program			
a = b * c ;			
if( a > c ){			
c = a * c - b;			
}			
0	lw	\$1, 24(\$fp)	
1	lw	\$2, 20(\$fp)	
2	mul	\$1, \$1, \$2	
3	sw	\$1, 28(\$fp)	
4	lw	\$1, 28(\$fp)	
5	lw	\$2, 20(\$fp)	
6	slt	\$1, \$2, \$1	
7	beqz	\$1, \$BB0_3	
8	j	\$B2	
9	\$B2:		
10	lw	\$1, 28(\$fp)	
11	lw	\$2, 20(\$fp)	
12	mul	\$1, \$1, \$2	
13	lw	\$2, 24(\$fp)	
14	subu	\$1, \$1, \$2	
15	j	\$B3	
16	sw	\$1, 20(\$fp)	

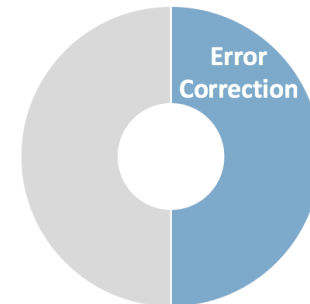
mem h0  
mem h1  
art h2  
mem h3  
mem h4  
mem h5  
art h6  
br h7  
br h8  
br h9  
mem h10  
mem h11  
art h12  
mem h13  
art h14  
br h15  
mem h16



# Stage 1: Coda Sketch Generation

- **Decoder**

- Generate Abstract Syntax Tree (AST)
- AST can be equivalently translated into its corresponding high level Program
- Advantages:
  - Prevent error propagation/ Preserve node dependency / capture PL grammar
  - Boundaries are more explicit (terminal nodes)
- Using Attention Mechanism



# Stage 2: Iterative Error Correction

- The sketch generated in Stage 1 may contain errors:
  - mispredicted tokens, missing lines, redundant lines

## Golden program

```
If( a > c ) {  
  a = b + c * a;  
  b = a - c;  
}
```

## Wrongly predicted

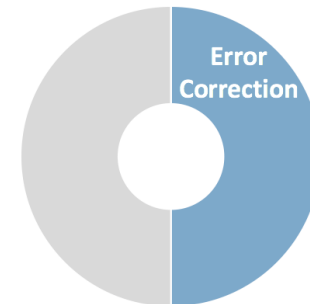
```
If( a > b ) {  
  a = b + c * a;  
  b = a - b;  
}
```

## Missing lines

```
If( a > c ) {  
  a = b + c * a;  
}
```

## Redundant lines

```
If( a > c ) {  
  a = b + c * a;  
  b = a;  
  b = a;  
}
```



# Stage 2: Iterative Error Correction

- Correct the error by iterative Error Predictor (EP)
  - Iterative rewriting!
  - Spot errors in the generated assembly codes
  - Fix errors and recompile
  - Repeat 10 times

# Experimental Setup

- Compiler configuration: Clang **-O0** (no code optimization)
- Benchmarks:
  - Synthetic programs:
    - **Karel library (Karel)** – only function calls
    - **Math library (Math)** – function calls with arguments
    - **Normal expressions (NE)** – (^,&,\*,-,<<, >>,|,% ...)
    - **Math library + Normal expressions (Math+NE)** – replaces the variables in NE with a return value of math function.
- Metrics:
  - Token Accuracy
  - Program Accuracy

# Result – Stage 1 Performance

- Token accuracy (%) across benchmarks

Benchmarks	Seq2Seq	Seq2Seq+Attn	Seq2AST+Attn	Inst2seq+Attn	Inst2AST+Attn
Karel <sub>S</sub>	51.61	97.13	99.81	98.83	<b>99.89</b>
Math <sub>S</sub>	23.12	94.85	99.12	96.20	<b>99.72</b>
NE <sub>S</sub>	18.72	87.36	90.45	88.48	<b>94.66</b>
(Math+NE) <sub>S</sub>	14.14	87.86	91.98	89.67	<b>97.90</b>
Karel <sub>L</sub>	33.54	94.42	98.02	98.12	<b>98.56</b>
Math <sub>L</sub>	11.32	91.94	96.63	93.16	<b>98.63</b>
NE <sub>L</sub>	11.02	81.80	85.92	85.97	<b>91.92</b>
(Math+NE) <sub>L</sub>	6.09	81.56	85.32	86.16	<b>93.20</b>

- Highest token accuracy across all benchmarks (96.8% on average) compared to baselines.
- 10.1% and 80.9% margin over a naive Seq2Seq model with and without attention.
- More tolerant to the growth of program length.



# Result – Stage 2 Performance

- Program accuracy (%)

BenchMarks	(i) Error Detection		(ii) Befor EC		After EC	
	<i>s2s</i> ,10	<i>i2a</i> ,10	<i>s2s</i>	<i>i2a</i>	<i>s2s</i>	<i>i2a</i>
Math <sub>S</sub>	91.4	94.2	40.1	64.8	91.2	<b>100.0</b>
NE <sub>S</sub>	83.5	88.7	6.6	12.2	53.0	<b>78.6</b>
(Math+NE) <sub>S</sub>	83.6	90.1	3.5	43.2	63.6	<b>89.2</b>
Math <sub>L</sub>	87.5	91.3	21.7	51.8	83.9	<b>99.5</b>
NE <sub>L</sub>	78.1	84.5	0.2	2.6	33.1	<b>56.4</b>
(Math+NE) <sub>L</sub>	80.2	85.3	0.1	4.9	38.3	<b>67.2</b>

Baseline

Ours

s2s = sequence-to-sequence with attention

i2a = instruction encoder to AST decoder with attention

# Summarization and Future Works

- Summary
  - Gives examples of scalable RL system
  - RL/ML can be used to learn heuristics for system
- Large Open Space Ahead
  - ML captures statistics regularity and leads to better solutions
  - Application to large-scale systems?
  - Theoretical Guarantees?



Thanks!