Device Placement Optimization with Reinforcement Learning

Azalia Mirhoseini, Anna Goldie, Hieu Pham, Benoit Steiner, Mohammad Norouzi, Naveen Kumar, Rasmus Larsen, Yuefeng Zhou, Quoc Le, Samy Bengio, and Jeff Dean



Why device placement?

Trend towards many-device training, bigger models, larger batch sizes



Google neural machine translation'16 (trained on 128 GPUs)





Sparsely gated mixture of experts'17 (batch size = 8192, >130 billion parameters, trained on 128 GPUs)

Training ImageNet in 1-hr'17 (batch size = 8192, trained on 256 GPUs)

Standard practice for device placement

- Often based on greedy heuristics
- Requires deep understanding of devices: nonlinear FLOPs, bandwidth, latency behavior
- Requires modeling parallelism and pipelining
- Does not generalize well

ML for device placement

• ML is repeatedly replacing rule based heuristics

ML for device placement

- ML is repeatedly replacing rule based heuristics
- We show how RL can be applied to device placement
 - Effective search across large state and action spaces
 - Automated learning from environment only based on reward function (e.g. runtime of a program)

ML for device placement

- ML is repeatedly replacing rule based heuristics
- We show how RL can be applied to device placement
 - Effective search across large state and action spaces
 - Automated learning from environment only based on reward function (e.g. runtime of a program)
- We are inspired by success of ML for learning to search

arXiv.org > cs > arXiv:1611.09940	arXiv.org > cs > arXiv:1611.01578	arXiv.org > stat > arXiv:1506.03134	
Computer Science > Artificial Intelligence	Computer Science > Learning	Statistics > Machine Learning	
Neural Combinatorial Optimization with Reinforcement Learning	Neural Architecture Search with Reinforcement Learning	Pointer Networks Oriol Vinyals, Meire Fortunato, Navdeep Jaitly	
Irwan Bello, Hieu Pham, Quoc V. Le, Mohammad Norouzi, Samy Bengio	Barret Zoph, Quoc V. Le		

Posing device placement as an RL problem



Posing device placement as an RL problem



Problem formulation

$$J(\theta) = \mathbf{E}_{\mathcal{P} \sim \pi(\mathcal{P}|\mathcal{G};\theta)} \left[R\left(\mathcal{P}\right) | \mathcal{G} \right]$$

 $J(\theta): expected runtime \\ \theta: trainable parameters of policy \\ R: runtime \\ \pi(\mathcal{P}|g;\theta): policy \\ g: input neural graph \\ \mathcal{P}: output placements \in \{1,2,..,num_ops\}^{num_devices}$

Policy architecture



Policy architecture



Training with **REINFORCE**

$$J(\theta) = \mathbf{E}_{\mathcal{P} \sim \pi(\mathcal{P}|\mathcal{G};\theta)} \left[R\left(\mathcal{P}\right) | \mathcal{G} \right]$$

 $\nabla_{\theta} J(\theta) = \mathbf{E}_{\mathcal{P} \sim \pi(\mathcal{P}|\mathcal{G};\theta)} \left[R\left(\mathcal{P}\right) \cdot \nabla_{\theta} \log p\left(\mathcal{P}|\mathcal{G};\theta\right) \right]$

Training with **REINFORCE**

$$J(\theta) = \mathbf{E}_{\mathcal{P} \sim \pi(\mathcal{P}|\mathcal{G};\theta)} \left[R\left(\mathcal{P}\right) | \mathcal{G} \right]$$

$$\nabla_{\theta} J(\theta) = \mathbf{E}_{\mathcal{P} \sim \pi(\mathcal{P}|\mathcal{G};\theta)} \left[R\left(\mathcal{P}\right) \cdot \nabla_{\theta} \log p\left(\mathcal{P}|\mathcal{G};\theta\right) \right]$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{K} \sum_{i=1}^{K} \left(R\left(\mathcal{P}_{i}\right) - B\right) \cdot \nabla_{\theta} \log p\left(\mathcal{P}_{i}|\mathcal{G};\theta\right)$$

Distributed training



Learned placement on Neural Machine Translation (NMT)



White represents CPU (Ixion Haswell 2300) Each other color represents a separate GPU (Nvidia Tesla K80) Searching over a space of 5²80 possible assignments

NMT end to end training



Mirhoseini, Pham et al., "Device Placement Optimization with Reinforcement Learning", ICML'17

Learned placement on Inception-V3



White represents CPU (Ixion Haswell 2300) Each other color represents a separate GPU (Nvidia Tesla K80) Searching over a space of 5⁸³ possible assignments

Inception-V3 end to end training



Profiling placement on NMT



Profiling placement on Inception-V3



Profiling placement on Inception-V3



Shortcomings of initially proposed model

- Seq2seq models cannot be unrolled for more than few hundred steps
- Most TensorFlow graphs contain tens of thousands of operations
- Manual grouping of operations hampers scalability



An end-to-end hierarchical placement model



Objective: Minimize expected runtime for predicted placement d

$$J(\theta_g, \theta_d) = \mathbf{E}_{\mathbf{P}(\mathbf{d}; \theta_{\mathbf{g}}, \theta_{\mathbf{d}})}[R_d] = \sum_{g \sim \pi_g} \sum_{d \sim \pi_d} p(g; \theta_g) p(d|g; \theta_d) R_d$$

 $J(\theta_{g}, \theta_{d})$: expected runtime θ_{g} : trainable parameters of Grouper θ_{d} : trainable parameters of Placer R_{d} : runtime for placement d

Objective: Minimize expected runtime for predicted placement d

$$J(\theta_g, \theta_d) = \mathbf{E}_{\mathbf{P}(\mathbf{d}; \theta_{\mathbf{g}}, \theta_{\mathbf{d}})}[R_d] = \sum_{g \sim \pi_g} \sum_{d \sim \pi_d} p(g; \theta_g) p(d|g; \theta_d) R_d$$

 $J(\theta_{g}, \theta_{d})$: expected runtime θ_{g} : trainable parameters of Grouper θ_{d} : trainable parameters of Placer R_{d} : runtime for placement d

$$J(\theta_g, \theta_d) = \mathbf{E}_{\mathbf{P}(\mathbf{d}; \theta_g, \theta_d)}[R_d] = \sum_{g \sim \pi_g} \sum_{d \sim \pi_d} p(g; \theta_g) p(d|g; \theta_d) R_d$$

Probability of predicted group assignment of

Probability of predicted group assignment operations

$$J(\theta_g, \theta_d) = \mathbf{E}_{\mathbf{P}(\mathbf{d}; \theta_{\mathbf{g}}, \theta_{\mathbf{d}})}[R_d] = \sum_{g \sim \pi_g} \sum_{d \sim \pi_d} p(g; \theta_g) p(d|g; \theta_d) R_d$$

Probability of predicted device placement conditioned on grouping results

Gradient update for Grouper

$$J(\theta_g, \theta_d) = \mathbf{E}_{\mathbf{P}(\mathbf{d}; \theta_{\mathbf{g}}, \theta_{\mathbf{d}})}[R_d] = \sum_{g \sim \pi_g} \sum_{d \sim \pi_d} p(g; \theta_g) p(d|g; \theta_d) R_d$$

$$\boxed{\nabla_{\theta g} J(\theta_g, \theta_d)} = \sum_{g \sim \pi_g} \nabla_{\theta g} p(g; \theta_g) \sum_{d \sim \pi_d} p(d|g; \theta_d) R_d$$

Derivative w.r.t. parameters of Grouper
$$\approx \frac{1}{m} \sum_{g_i \sim \pi_g}^{1 \leq i \leq m} \nabla_{\theta g} \log p(g_i; \theta_g) \cdot \frac{1}{k} (\sum_{d_j \sim \pi_d}^{1 \leq j \leq k} R_{d_j})$$

Gradient update for Placer

$$J(\theta_g, \theta_d) = \mathbf{E}_{\mathbf{P}(\mathbf{d}; \theta_{\mathbf{g}}, \theta_{\mathbf{d}})}[R_d] = \sum_{g \sim \pi_g} \sum_{d \sim \pi_d} p(g; \theta_g) p(d|g; \theta_d) R_d$$

$$\nabla_{\theta d} J(\theta_g, \theta_d) = \sum_{d \sim \pi_d} \sum_{g \sim \pi_g} p(g; \theta_g) \nabla_{\theta d} p(d|g; \theta_d) R_d$$
Derivative w.r.t. parameters of Placer
$$\approx \frac{1}{k} \sum_{d_j \sim \pi_d}^{1 \leq j \leq k} \frac{1}{m} (\sum_{g_i \sim \pi_g}^{1 \leq i \leq m} \nabla_{\theta d} \log p(d_j | g_i; \theta_d) R_{d_j})$$

Learned placements on NMT



White represents CPU (Ixion Haswell 2300) Each other color represents a separate GPU (Nvidia Tesla K80) Searching over a space of ~5^40000 possible assignments

Results (runtime in seconds)

Tasks	CPU	GPU	#GPUs	Human	Scotch	MinCut	Hierarchical	Runtime
	Only	Only		Expert			Planner	Reduction
Inception-V3	0.61	0.15	2	0.15	0.93	0.82	0.13	16.3%
ResNet	-	1.18	2	1.18	6.27	2.92	1.18	0%
RNNLM	6.89	1.57	2	1.57	5.62	5.21	1.57	0%
NMT (2-layer)	6.46	OOM	2	2.13	3.21	5.34	0.84	60.6%
NMT (4-layer)	10.68	OOM	4	3.64	11.18	11.63	1.69	53.7%
NMT (8-layer)	11.52	OOM	8	3.88	17.85	19.01	4.07	-4.9%

Results (runtime in seconds)

Tasks	CPU	GPU	#GPUs	Human	Scotch	MinCut	Hierarchical	Runtime
	Only	Only		Expert			Planner	Reduction
Inception-V3	0.61	0.15	2	0.15	0.93	0.82	0.13	16.3%
ResNet	-	1.18	2	1.18	6.27	2.92	1.18	0%
RNNLM	6.89	1.57	2	1.57	5.62	5.21	1.57	0%
NMT (2-layer)	6.46	OOM	2	2.13	3.21	5.34	0.84	60.6%
NMT (4-layer)	10.68	OOM	4	3.64	11.18	11.63	1.69	53.7%
NMT (8-layer)	11.52	OOM	8	3.88	17.85	19.01	4.07	-4.9%

Results (runtime in seconds)

Tasks	CPU	GPU	#GPUs	Human	Scotch	MinCut	Hierarchical	Runtime
	Only	Only		Expert			Planner	Reduction
Inception-V3	0.61	0.15	2	0.15	0.93	0.82	0.13	16.3%
ResNet	-	1.18	2	1.18	6.27	2.92	1.18	0%
RNNLM	6.89	1.57	2	1.57	5.62	5.21	1.57	0%
NMT (2-layer)	6.46	OOM	2	2.13	3.21	5.34	0.84	60.6%
NMT (4-layer)	10.68	OOM	4	3.64	11.18	11.63	1.69	53.7%
NMT (8-layer)	11.52	OOM	8	3.88	17.85	19.01	4.07	-4.9%

Summary

- Pose device placement as an RL problem
- Use policy gradient to learn placements
- Policy finds non-trivial assignment of operations to devices that outperform heuristic approaches
- Profiling of results show policy learns implicit trade-offs between computational and communication capabilities of underlying devices