



# PASHA: Efficient HPO and NAS with Progressive Resource Allocation

Ondrej Bohdal<sup>1\*</sup>, Lukas Balles<sup>2</sup>, Martin Wistuba<sup>2</sup>, Beyza Ermis<sup>3\*</sup>,  
Cedric Archambeau<sup>2</sup>, Giovanni Zappella<sup>2</sup>



THE UNIVERSITY of EDINBURGH  
**informatics**

<sup>1</sup>



| science

<sup>2</sup>



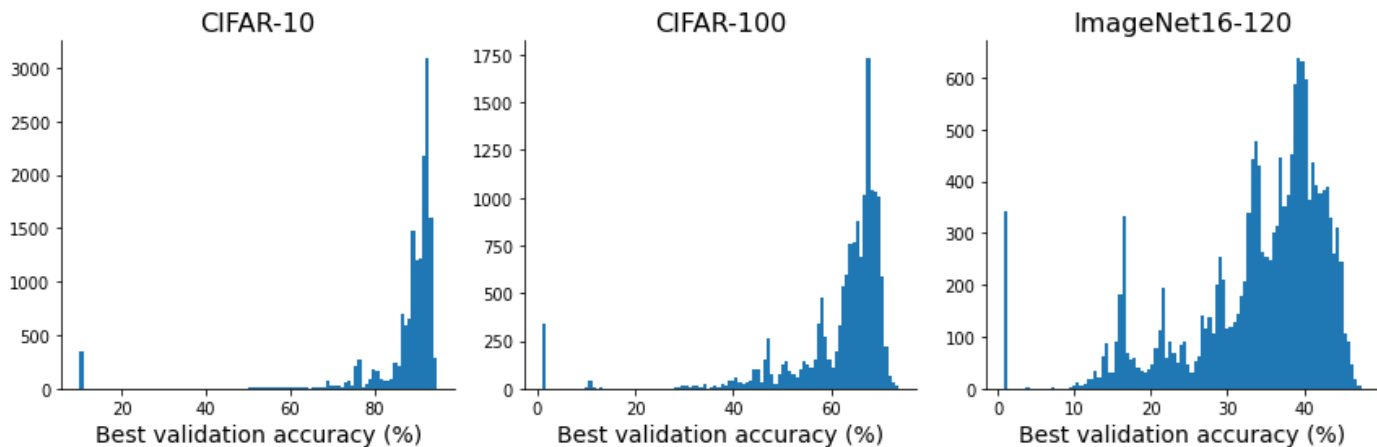
Cohere For AI

<sup>3</sup>

\* work done while at AWS, Berlin

# HPO and NAS

- Selection of hyperparameters and neural network architecture has a large impact on performance of the model
- Example: various architectures in NASBench201 have significantly different performances



# Challenges with Large-Scale HPO and NAS

- Costs make large-scale HPO difficult and often unviable
  - Evaluating 50 configurations for a 340-million-parameter BERT model (Devlin et al., NAACL'19) on the 15GB Wikipedia and Book corpora would cost around \$500,000
- Multi-fidelity methods such as ASHA (Li et al., MLSys'20) require specifying max amount of resources
- How to specify the max amount of resources?
  - Usually overestimated to guarantee convergence
  - Excellent configurations could be found using far fewer resources
- Insight: ranking of configurations is relatively stable after initial part of training
  - Learning curves rarely cross in later stages of training (excluding noise)

# PASHA

- Variation on ASHA (Li et al., MLSys'20) - asynchronous successive halving
- Idea: dynamically increase max resources depending on if ranking of configurations is stable
  - Start with a small initial amount of maximum resources
  - Progressively increase them if the ranking of the configurations in the top two rungs (rounds of promotion) has not stabilized
  - Due to stochasticity, some benevolence in rankings needed → use soft ranking
- Particularly useful for HPO on massive datasets

# Illustration

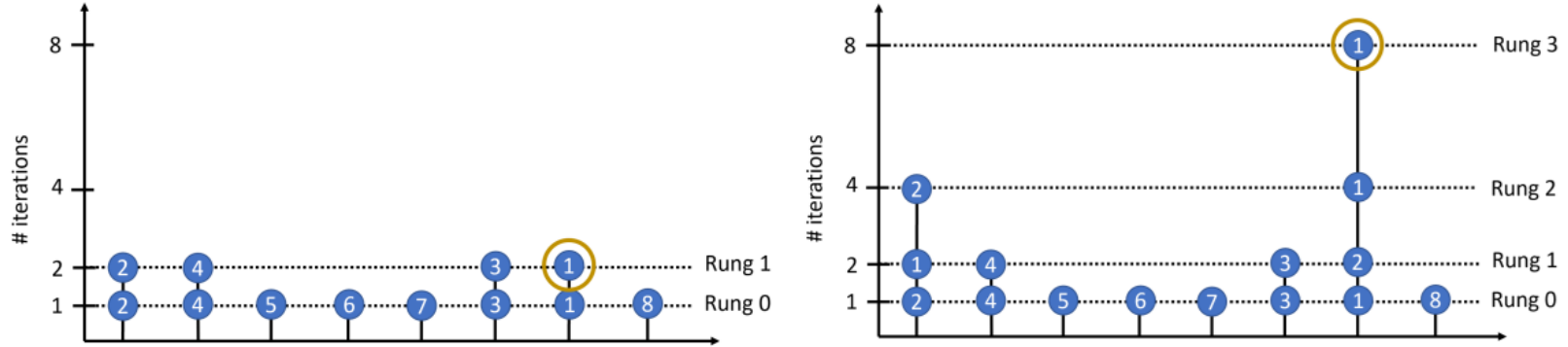


Figure 1: Illustration of how PASHA stops early if the ranking of configurations has stabilized. Left: the ranking of the configurations (displayed inside the circles) has stabilized, so we can select the best configuration and stop the search. Right: the ranking has not stabilized, so we continue.

# Soft Ranking

- Configurations are equivalent if their performance difference is less than  $\epsilon$ 
  - Estimate  $\epsilon$  based on noise in rankings across epochs
  - Intuition: configurations that repeatedly swap their rankings are similar

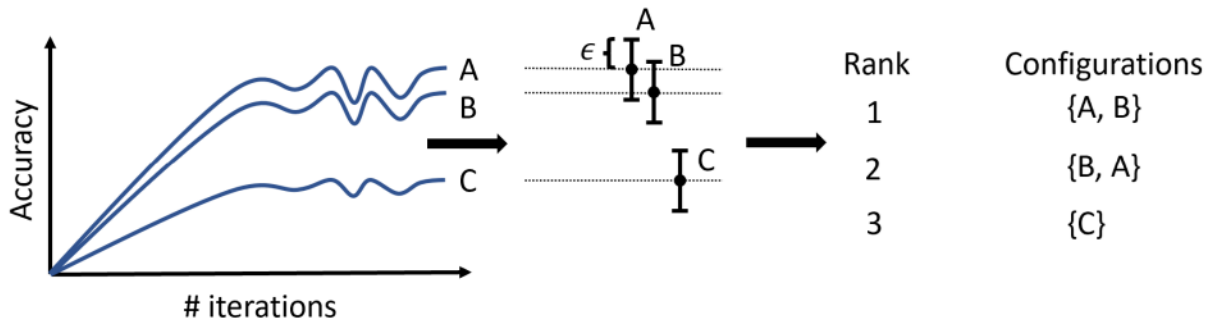


Figure 2: Illustration of soft ranking. There are three lists with the first two containing two items because the scores of the two configurations are closer to each other than  $\epsilon$ .

# Results - NAS

Table 1: NASBench201 results. PASHA leads to large improvements in runtime, while achieving similar accuracy as ASHA.

Dataset	Approach	Accuracy (%)	Runtime	Speedup factor	Max resources
CIFAR-10	ASHA	$93.85 \pm 0.25$	$3.0\text{h} \pm 0.6\text{h}$	1.0x	$200.0 \pm 0.0$
	PASHA	$93.57 \pm 0.75$	$1.3\text{h} \pm 0.6\text{h}$	2.3x	$36.1 \pm 50.0$
	One-epoch baseline	$93.30 \pm 0.61$	$0.3\text{h} \pm 0.0\text{h}$	8.5x	$1.0 \pm 0.0$
	Random baseline	$72.88 \pm 19.20$	$0.0\text{h} \pm 0.0\text{h}$	N/A	$0.0 \pm 0.0$
CIFAR-100	ASHA	$71.69 \pm 1.05$	$3.2\text{h} \pm 0.9\text{h}$	1.0x	$200.0 \pm 0.0$
	PASHA	$71.84 \pm 1.41$	$0.9\text{h} \pm 0.4\text{h}$	3.4x	$20.5 \pm 48.3$
	One-epoch baseline	$65.57 \pm 5.53$	$0.3\text{h} \pm 0.0\text{h}$	9.2x	$1.0 \pm 0.0$
	Random baseline	$42.83 \pm 18.20$	$0.0\text{h} \pm 0.0\text{h}$	N/A	$0.0 \pm 0.0$
ImageNet16-120	ASHA	$45.63 \pm 0.81$	$8.8\text{h} \pm 2.2\text{h}$	1.0x	$200.0 \pm 0.0$
	PASHA	$45.13 \pm 1.51$	$2.9\text{h} \pm 1.7\text{h}$	3.1x	$21.3 \pm 48.1$
	One-epoch baseline	$41.42 \pm 4.98$	$1.0\text{h} \pm 0.0\text{h}$	8.8x	$1.0 \pm 0.0$
	Random baseline	$20.75 \pm 9.97$	$0.0\text{h} \pm 0.0\text{h}$	N/A	$0.0 \pm 0.0$

# Combination with Bayesian Optimization

Table 2: NASBench201 results for ASHA with Bayesian Optimization searcher – MOBSTER (Klein et al., 2020) and similarly extended version of PASHA. The results show PASHA can be successfully combined with a smarter configuration selection strategy.

Dataset	Approach	Accuracy (%)	Runtime	Speedup factor	Max resources
CIFAR-10	MOBSTER	$94.21 \pm 0.28$	$5.0\text{h} \pm 1.1\text{h}$	1.0x	$200.0 \pm 0.0$
	PASHA BO	$94.00 \pm 0.20$	$2.6\text{h} \pm 1.8\text{h}$	2.0x	$70.7 \pm 81.6$
CIFAR-100	MOBSTER	$72.79 \pm 0.68$	$5.7\text{h} \pm 1.4\text{h}$	1.0x	$200.0 \pm 0.0$
	PASHA BO	$72.16 \pm 1.07$	$1.6\text{h} \pm 0.5\text{h}$	3.7x	$13.0 \pm 8.7$
ImageNet16-120	MOBSTER	$46.21 \pm 0.70$	$15.1\text{h} \pm 4.0\text{h}$	1.0x	$200.0 \pm 0.0$
	PASHA BO	$45.36 \pm 1.06$	$3.9\text{h} \pm 1.2\text{h}$	3.9x	$11.8 \pm 7.9$



# Results - HPO on Large Datasets

Table 3: Results of the HPO experiments on WMT and ImageNet tasks from the PD1 benchmark. Mean and std of the best validation accuracy (or its equivalent as given in the PD1 benchmark).

Dataset	Approach	Accuracy (%)	Runtime	Speedup factor	Max resources
WMT	ASHA	$62.72 \pm 1.41$	$43.7h \pm 37.2h$	1.0x	$1357.4 \pm 80.4$
	PASHA	$62.04 \pm 2.05$	$2.8h \pm 0.6h$	15.5x	$37.8 \pm 21.6$
	One-epoch baseline	$62.36 \pm 1.40$	$0.6h \pm 0.0h$	67.3x	$1.0 \pm 0.0$
	Random baseline	$33.93 \pm 21.96$	$0.0h \pm 0.0h$	N/A	$0.0 \pm 0.0$
ImageNet	ASHA	$75.10 \pm 2.03$	$7.3h \pm 1.2h$	1.0x	$251.0 \pm 0.0$
	PASHA	$73.37 \pm 2.71$	$3.8h \pm 1.0h$	1.9x	$45.0 \pm 30.1$
	One-epoch baseline	$63.40 \pm 9.91$	$1.1h \pm 0.0h$	6.7x	$1.0 \pm 0.0$
	Random baseline	$36.94 \pm 31.05$	$0.0h \pm 0.0h$	N/A	$0.0 \pm 0.0$

# Summary

- PASHA dynamically selects the amount of maximum resources
- Significant speedup of HPO and NAS without sacrificing the performance
  - Especially on large datasets
- Can be combined with Bayesian Optimization search strategies
- PASHA is available within Syne Tune HPO framework:  
<https://github.com/aws-labs/syne-tune>
- Tutorial for PASHA is also available: <https://syne-tune.readthedocs.io/en/latest/tutorials/pasha/pasha.html>