# Results of the Poll
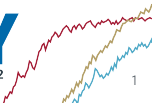Note: Answers labelled "academic" are mostly PhD students

**HITY**
**NeurIPS Workshop 2022**

# Which optimizer (i.e. update rule) do you try out first to train a neural network?



Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

Adam/AdamW

80%

84% 76%

0% 0% 1%
0% 0%
5% 2%
0% 4% 2% 1%
5% 2%
7%
7%

4% — other

3% — Shampoo

2% — Nesterov Momentum (aka. Nesterov Accelerated Momentum, NAG)

AdaGrad

Heavy Ball Momentum (aka. Polyak Momentum, vanilla Momentum)

Vanilla SGD (without Momentum)
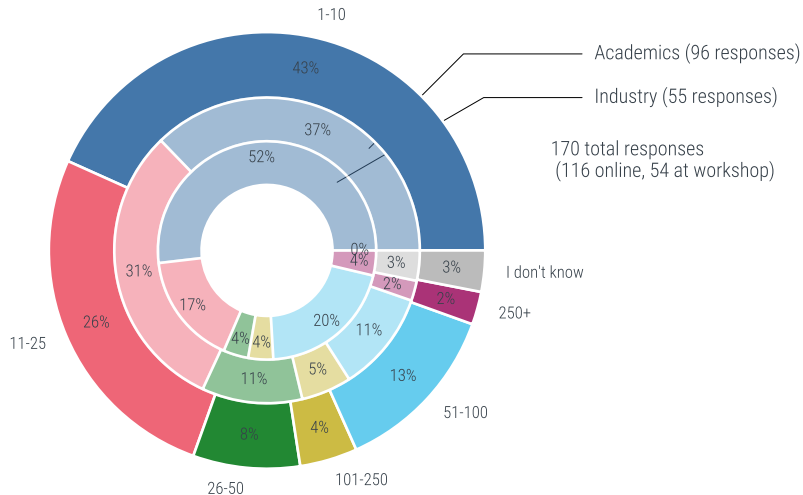
HITY

NeurIPS Workshop 2022

# Which optimizer (i.e. update rule) do you try out second to train a neural network? This is not a duplicate question! This question is about your second choice.



Heavy Ball Momentum (aka. Polyak Momentum, vanilla Momentum)

Nesterov Momentum (aka. Nesterov Accelerated Momentum, NAG)

Adam/AdamW

Vanilla SGD (without Momentum)

RMSProp

LARS/LAMB

AdaGrad

AdaDelta

Shampoo

RAdam

other

I don't need a second one

Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

HITY
NeurIPS Workshop 2022

3

How do you tune your hyperparameters?

Academics (96 responses)
Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

Manual tuning — 47%, 53%, 42%
Grid search — 29%, 28%, 11%
Random search — 13%, 12%, 2%
No tuning — 1%, 0%
Bayesian optimization — 9%, 5%, 16%
other — 2%, 1%, 0%, 0%

HITY
NeurIPS Workshop 2022

# How many tuning trials do you use? Just use your best estimate.



Academics (96 responses)
Industry (55 responses)
170 total responses
(116 online, 54 at workshop)

1-10: 43%, 37%, 52%
11-25: 31%, 26%, 17%
26-50: 4%, 4%, 11%, 8%
101-250: 4%, 5%, 11%, 13%
51-100: 20%, 2%, 2%
250+: 3%, 3%, 0%, 4%
I don't know

HITY
NeurIPS Workshop 2022

5

# How many hyperparameters do you tune?

Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

2-3 · 44% · 43% · 47%

1 · 6% · 2% · 2%

other · 2%

5+ · 18% · 17% · 17%

4-5 · 30% · 36% · 20%

8% · 7%

HITY
NeurIPS Workshop 2022

**Which hyperparameters do you tune?**

NeurIPS Workshop 2022

# Do you determine your batch size by memory, i.e. such that one batch fits exactly in memory?



Yes
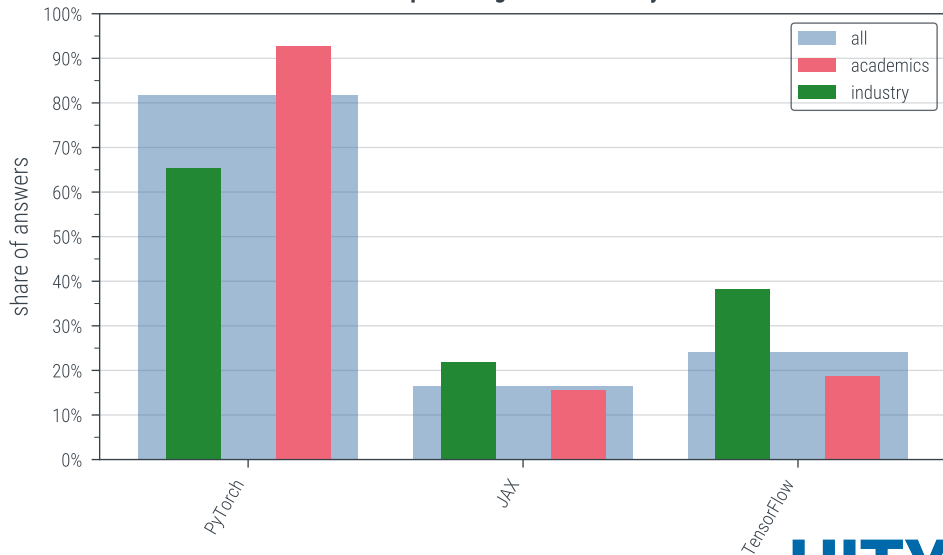
59%

52%

65%

Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

0% 0% 1%
0% 0%

2%

other

33%

45%

39%

No

HITY
NeurIPS Workshop 2022

8

# Which learning rate schedule do you try out first?



No schedule/constant learning rate — 33%, 32%, 28%

Linear decay — 20%, 16%, 28%

Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

other — 5%

Cyclical schedule — 3%

Cosine decay — 23%, 26%, 22%

Step decay — 16%, 19%, 15%

0%, 0%, 0%

4%, 3%

0%, 1%

0%, 0%, 0%

0%, 0%, 0%

HITY
NeurIPS Workshop 2022

# Do you use a learning rate warmup?



- Academics (96 responses)
- Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

Yes — 43%, 34%
No — 55%, 65%, 40%, 60%
other — 0%, 0%, 1%

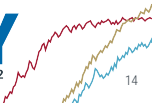# Which deep learning framework do you use?

- ▶ Marginal likelihood
- ▶ OneCycle scheduler, gradient checkpointing
- ▶ Genetic Algorithm for Hyperparameters
- ▶ W&B
- ▶ avoid batches that lead to nan/inf losses
- ▶ One cycle, low fidelity training, sgd with restarts
- ▶ Use proximal optimization for regularizers
- ▶ The Composer library for PyTorch (MosaicML)
- ▶ using line search to choose the maximum learning rate possible at each step
- ▶ Normalized updates
- ▶ We recently published a survey of our tricks: `https://arxiv.org/abs/2209.05310`
- ▶ Use Distributed Shampoo, Normformer, GLU
- ▶ Weight averaging

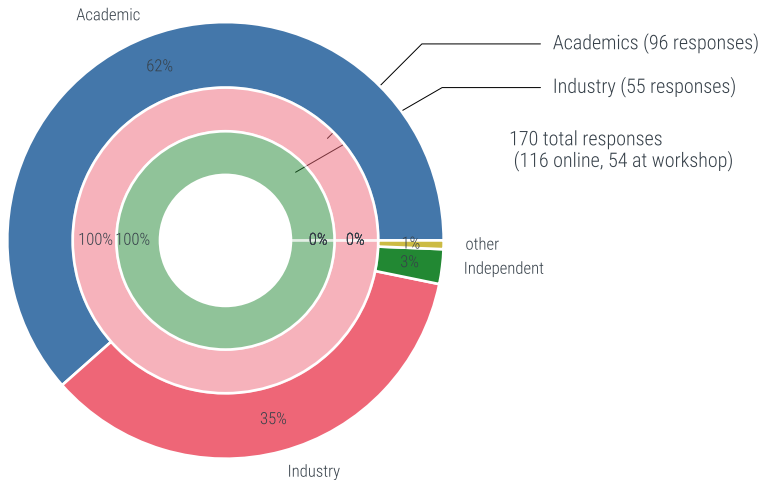# What are your favorite training tricks for neural networks that more people should know? Slide 2 of 2

► FreezeOut

► L2 regularised features

► try a different epsilon value!

► Use sample prioritization techniques, and perform calibration during training.

► free nats (ignoring losses that are smaller than some threshold)

► Check consistency of hyperparameter performance over multiple seeds

► Lowering the learning rate!

► In some contexts, normalizing data as pre-processing step works a lot better than batch or layer norm. Better than batch or layer norm.

► Mixed precision training

► label smoothing

► Train with a small subset

► Ablation spreadsheets and oncalls. :)

► Cyclic and one cycle LR

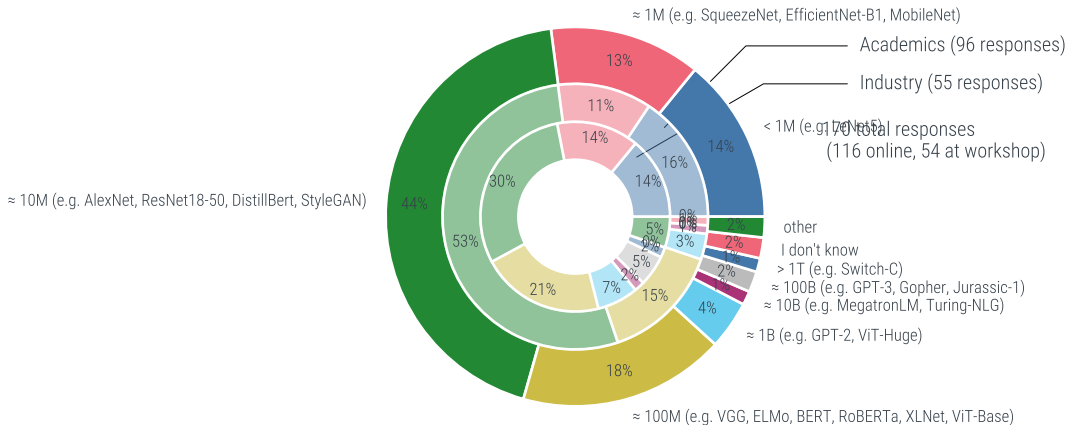► Label smoothing (easy & works well), looking at outputs for failure modes

– Backup –

**HITY**
NeurIPS Workshop 2022

# What is your background?



Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

How many weights do you usually have in your neural network architectures?
Please estimate the order of magnitude.

≈ 1M (e.g. SqueezeNet, EfficientNet-B1, MobileNet)

Academics (96 responses)

Industry (55 responses)

< 1M (e.g. LeNet-5) 170 total responses
(116 online, 54 at workshop)

≈ 10M (e.g. AlexNet, ResNet18-50, DistillBert, StyleGAN)

other

I don't know

> 1T (e.g. Switch-C)

≈ 100B (e.g. GPT-3, Gopher, Jurassic-1)

≈ 10B (e.g. MegatronLM, Turing-NLG)

≈ 1B (e.g. GPT-2, ViT-Huge)

≈ 100M (e.g. VGG, ELMo, BERT, RoBERTa, XLNet, ViT-Base)

# How many layers do you usually have in your neural network architectures?



Less than 10 — Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)

other
More than 100

51-100

11-50

How many training samples are typically in your data sets? Please estimate the order of magnitude.

≈ 10k (e.g. MNIST, CIFAR-10, IMDB, Omniglot)

Academics (96 responses)

Industry (55 responses)

170 total responses

≈ 1,000M (i.e. billions of samples)

≈ 100 (e.g. iris)

other

> 100M

≈ 10M (e.g. MovieLens 20M, WMT14 fr-en, LM1B)

≈ 1M (e.g. ImageNet2012, WMT14 de-en, Open Images v4)

≈ 100k (e.g. LibriSpeech, COCO 2014, OGBG-molpcba, SQuADv2)

HITY
NeurIPS Workshop 2022

How many GPUs (or similar accelerators) do you usually use to run a single training run? Please select the closest match. Note, this question does not ask how many GPUs you have available to perform runs in parallel, but how many GPUs you use to train a single model.

Academics (96 responses)

Industry (55 responses)

170 total responses
(116 online, 54 at workshop)