# UNREASONABLE EFFECTIVENESS OF NOISE

# TRAINING VIA HYPERPARAMETERS PATH OPTIMIZATION

**Mykola Maksymenko**, R&D Director, **SoftServe**

# ARXIV:1909.04013 (2019)

WITH

VLAD PUSHKARIOV, TECHNION
YONATHAN EFRONI, TECHNION
MACIEJ KOCH-JANUSZ, ETH ZURICH

- TRAINING IN A JOINT WEIGHTS-HYPERPARAMETER SPACE
- EFFICIENT SAMPLING OF THE HYPERPARAMETERS
- NON-STATIONARY SCHEDULING PROTOCOL
- FASTER AND BETTER TRAINED MODELS

# GOAL: OPTIMAL, SIMPLER, UNIVERSAL
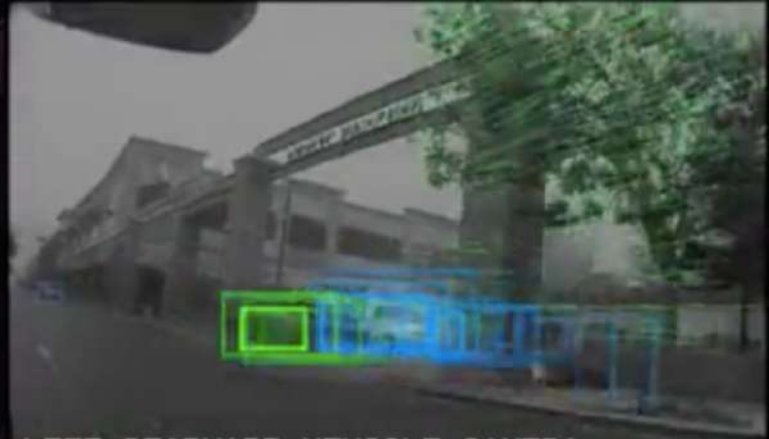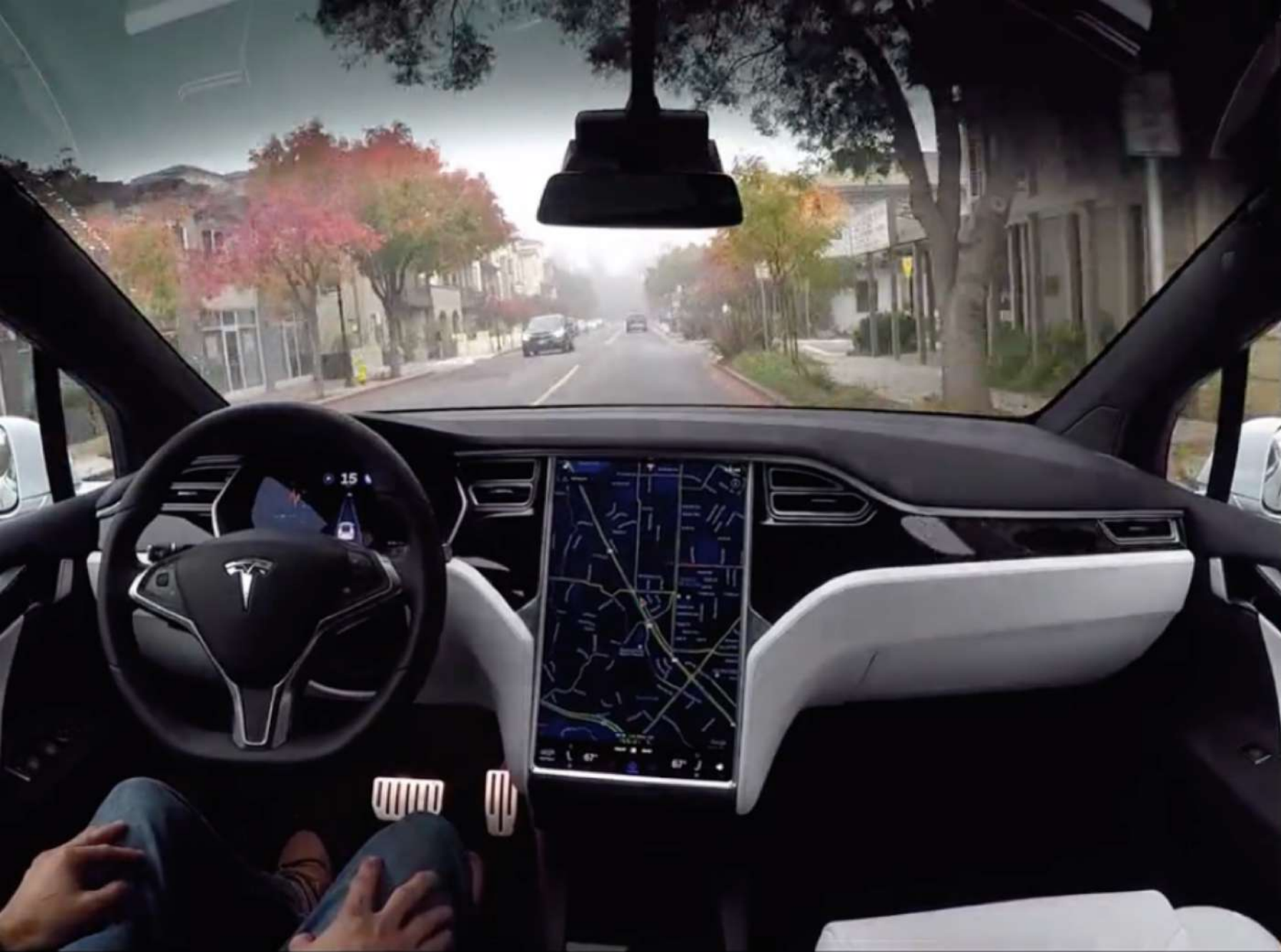
PHYSICS → MACHINE LEARNING

**Formalism** (energy, noise, etc)
Toy **models**
Concepts and **Ideas** (e.g. spin glass, tensor networks)
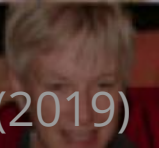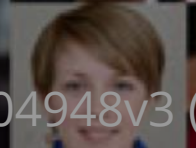etc...

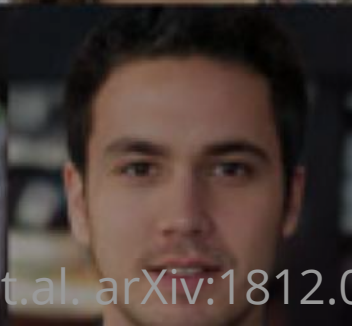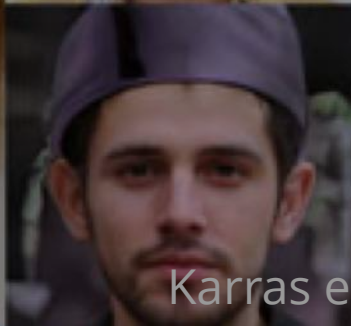softserve

LEFT REARWARD VEHICLE CAMERA

MEDIUM RANGE VEHICLE CAMERA

RIGHT REARWARD VEHICLE CAMERA

MOTION FLOW    LANE LINES    LANE LINES    ROAD FLOW    IN-PATH OBJECTS    ROAD LIGHTS    OBJECTS    ROAD SIGNS

# MECHANICS OF NEURAL NETWORKS

$$a_j^{(1)} = \sigma\left(\sum_i w_{ij}^{(1)} x_i + b_j^{(1)}\right)$$

$x_1$

$x_2$

$x_n$

$$a_j^{(f)} = \sigma\left(\sum_i w_{ij}^{(f)} a_i^{(1)} + b_j^{(f)}\right)$$
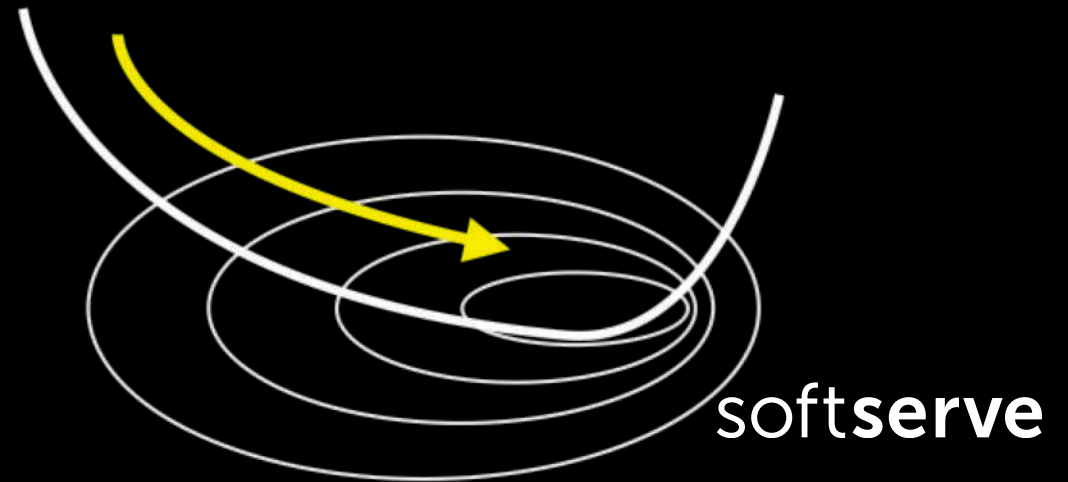
$$\text{Cost} = \frac{1}{2n} \sum_x \|a^f(x) - y_{\text{true}}(x)\|^2$$
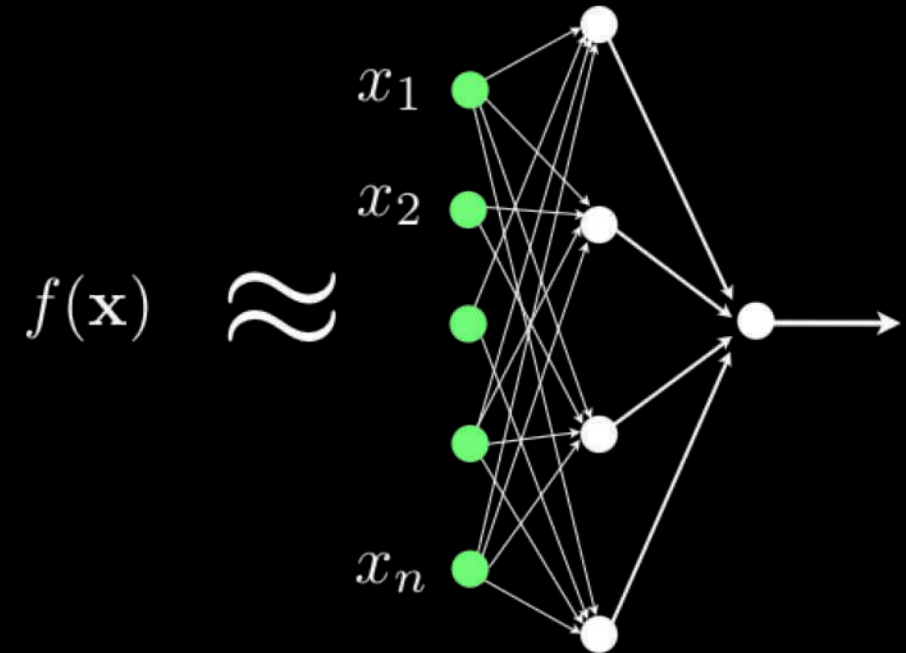
$$w_{ij_{t+1}} = w_{ij_t} - \gamma \nabla_{w_{ij}} Cost$$

softserve

# SHOULD WORK WELL IN THEORY...

Kolmogorov-Arnold representation theorem (1956-1963)

Universal approximation theorem, Cybenko (1989)

$$f(x_1, \ldots x_n) =$$

$$\sum_{j=1}^{2n+1} g_j \left( \sum_{i=1}^{n} \phi_{ij}(x_i) \right)$$

$$f(\mathbf{x}) \approx$$



softserve

40+ YEARS TO ADOPT IN PRACTICE

# DEEP ARCHITECTURE + BAG OF TRICKS

DEEP ARCHITECTURE
+ BAG OF TRICKS

BATCH GRADIENT
DROUPOUT
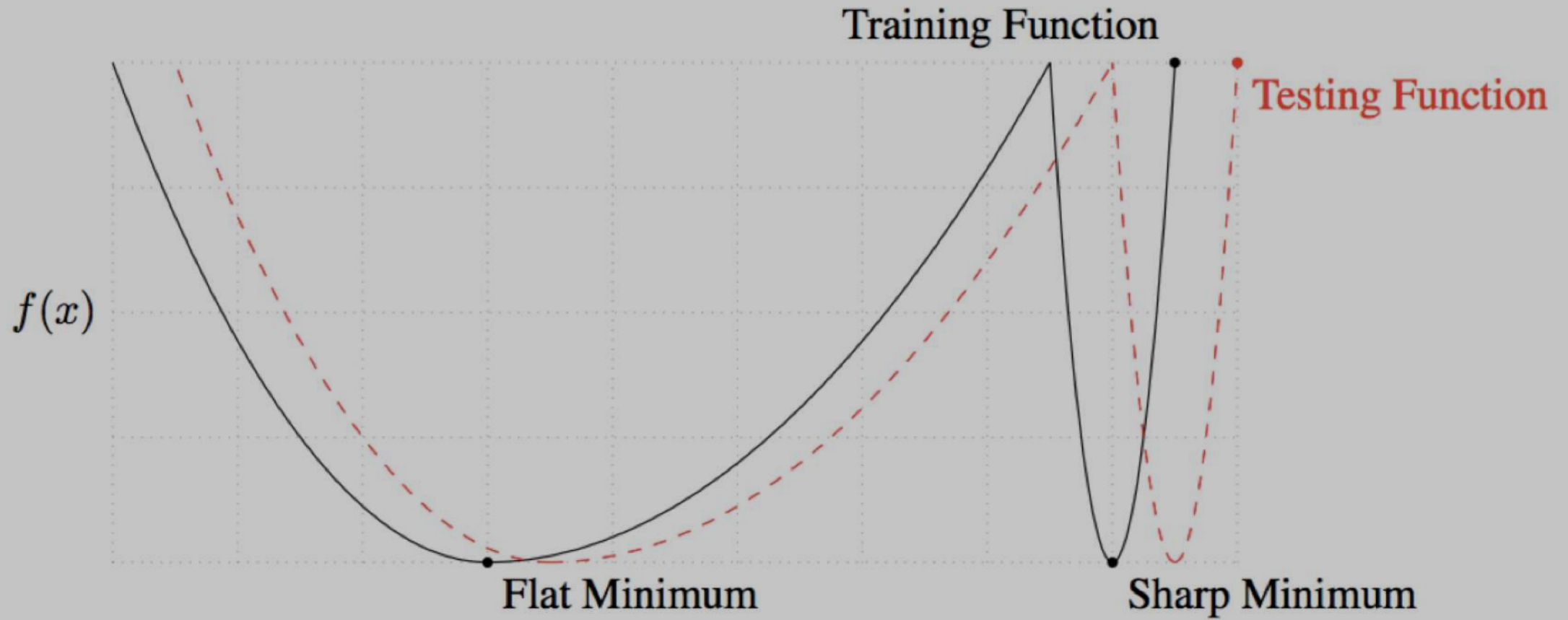EARLY STOPPING
ETC...

# DEEP ARCHITECTURE
# + BAG OF ~~TRICKS~~ NOISE
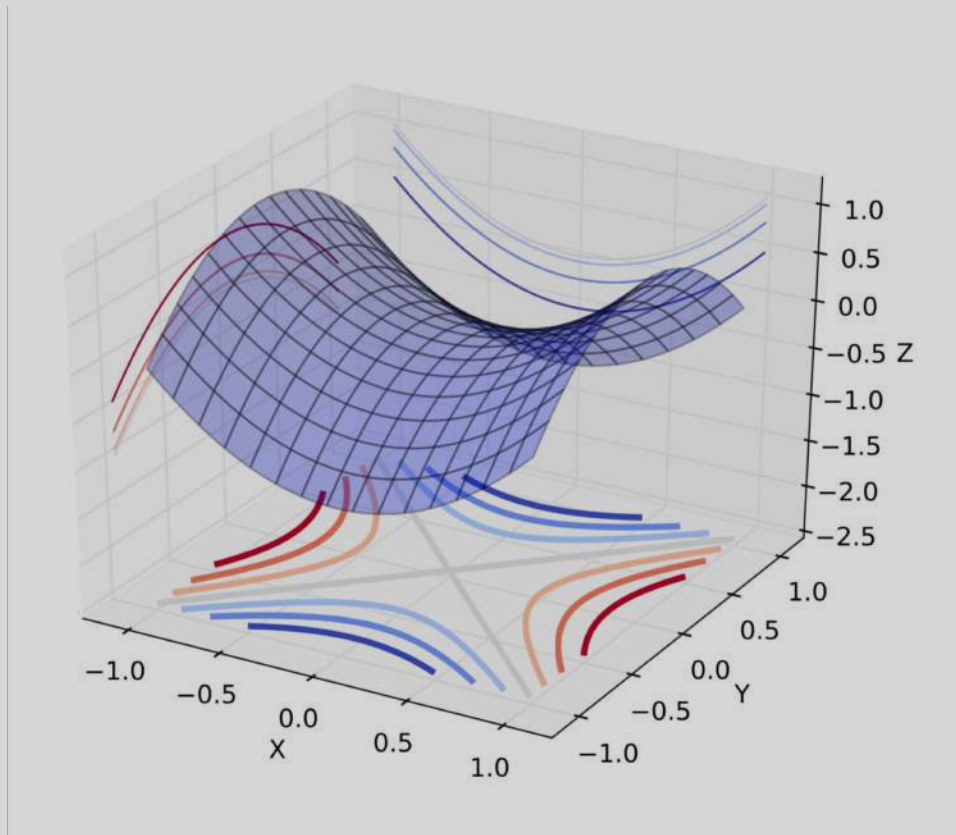
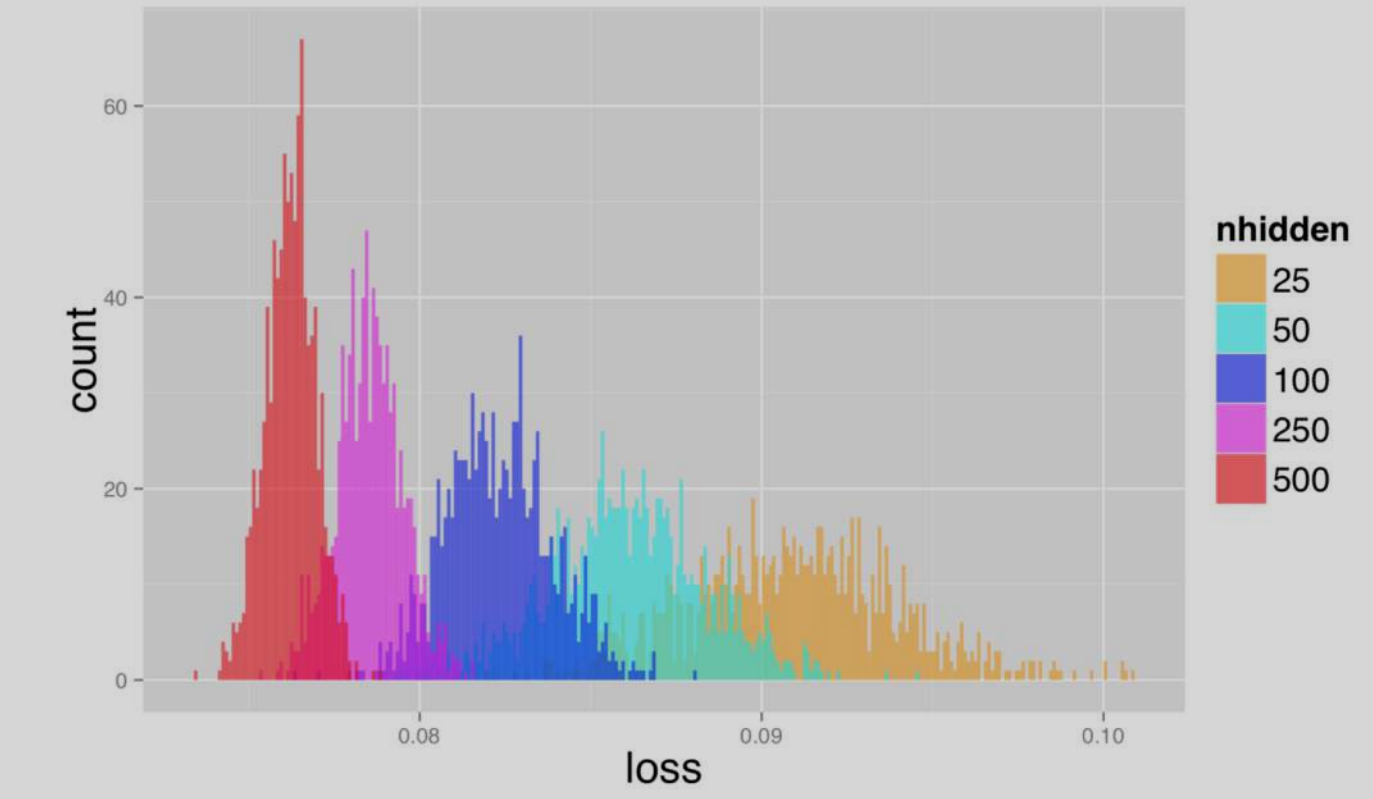# BATCH GRADIENT
# DROUPOUT
# EARLY STOPPING
# ETC...

# COMMON WISDOM: BATCH NOISE HELPS TO AVOID "BAD" MINIMA



Training Function

Testing Function

$f(x)$

Flat Minimum

Sharp Minimum

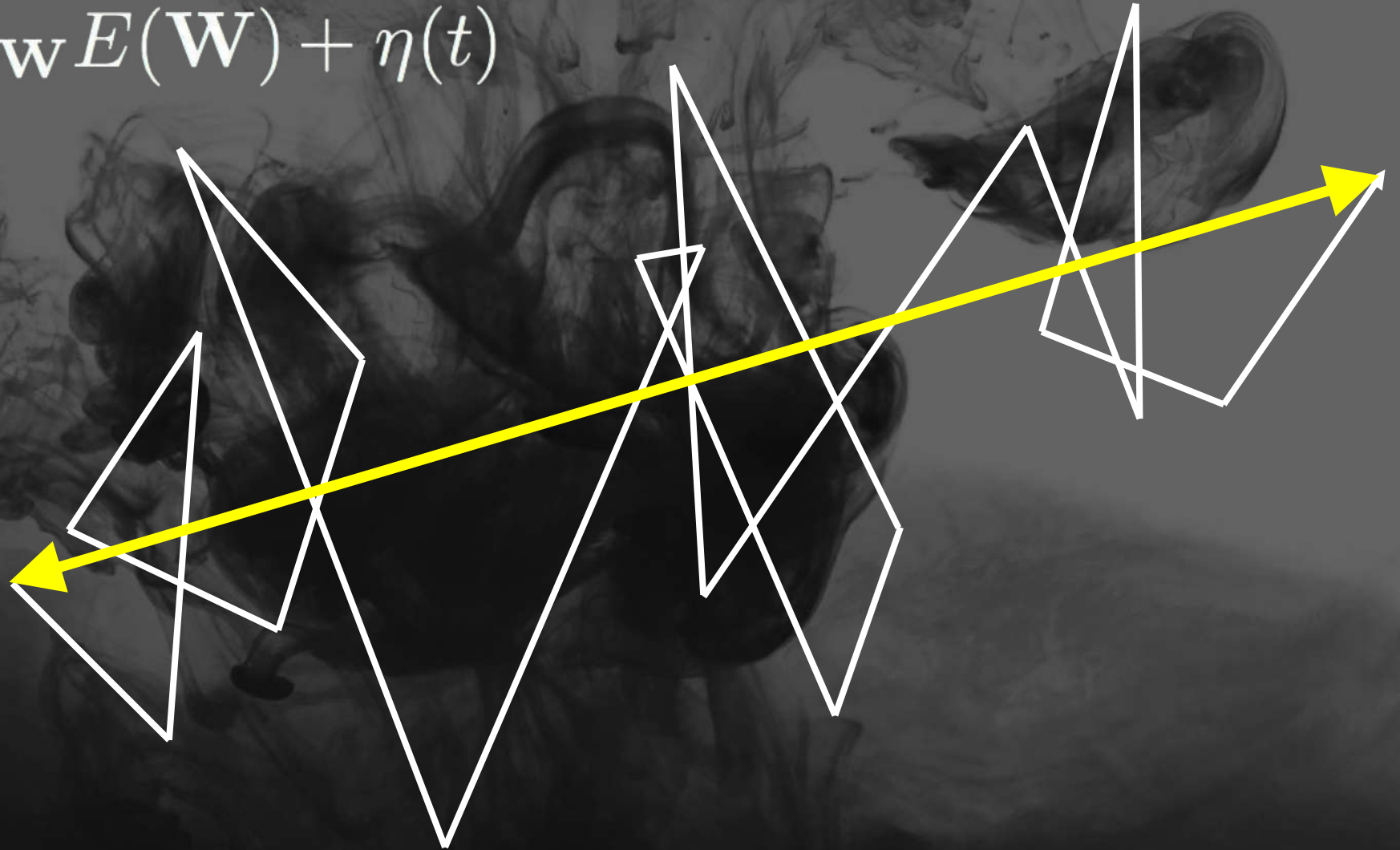# HOW LIKELY ARE THE "BAD" MINIMA?



arXiv:1406.2572v1 (2014)

arXiv:1412.0233v3 (2015)

**VERY UNLIKELY!**

# NOISE SPEEDS UP THE DIFFUSION!

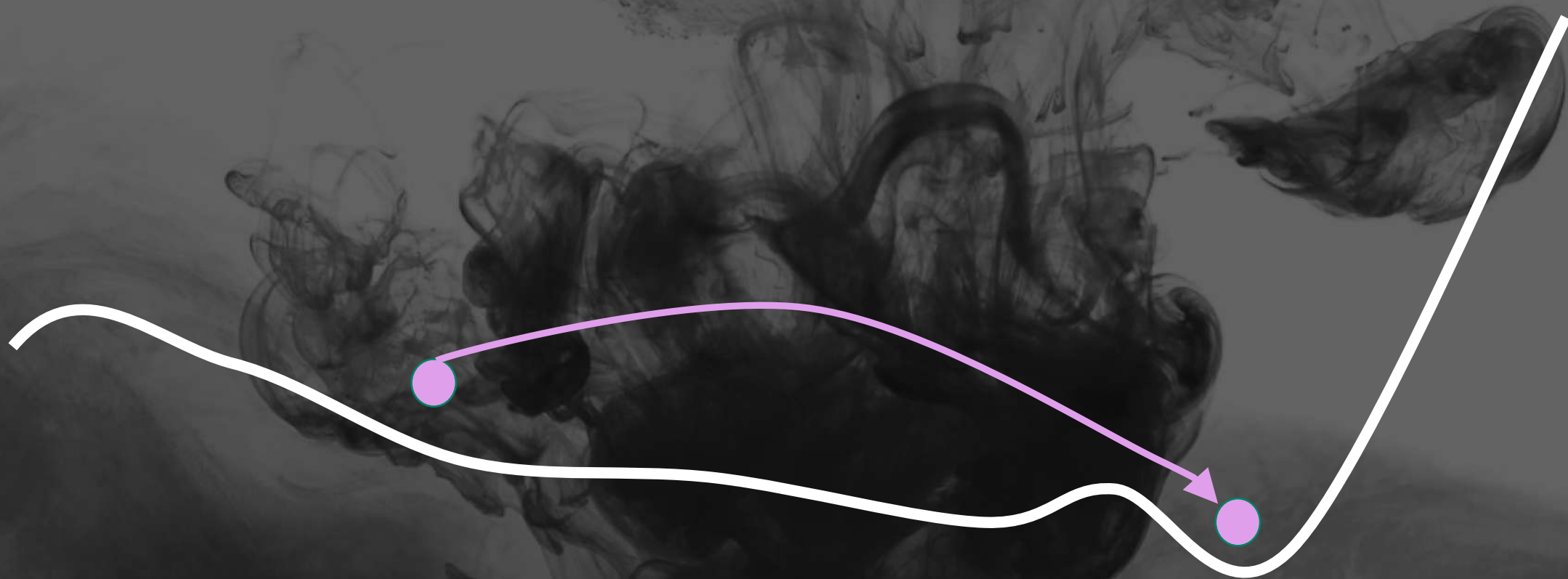$$\frac{\partial \mathbf{W}}{\partial t} = -\nabla_{\mathbf{W}} E(\mathbf{W}) + \eta(t)$$
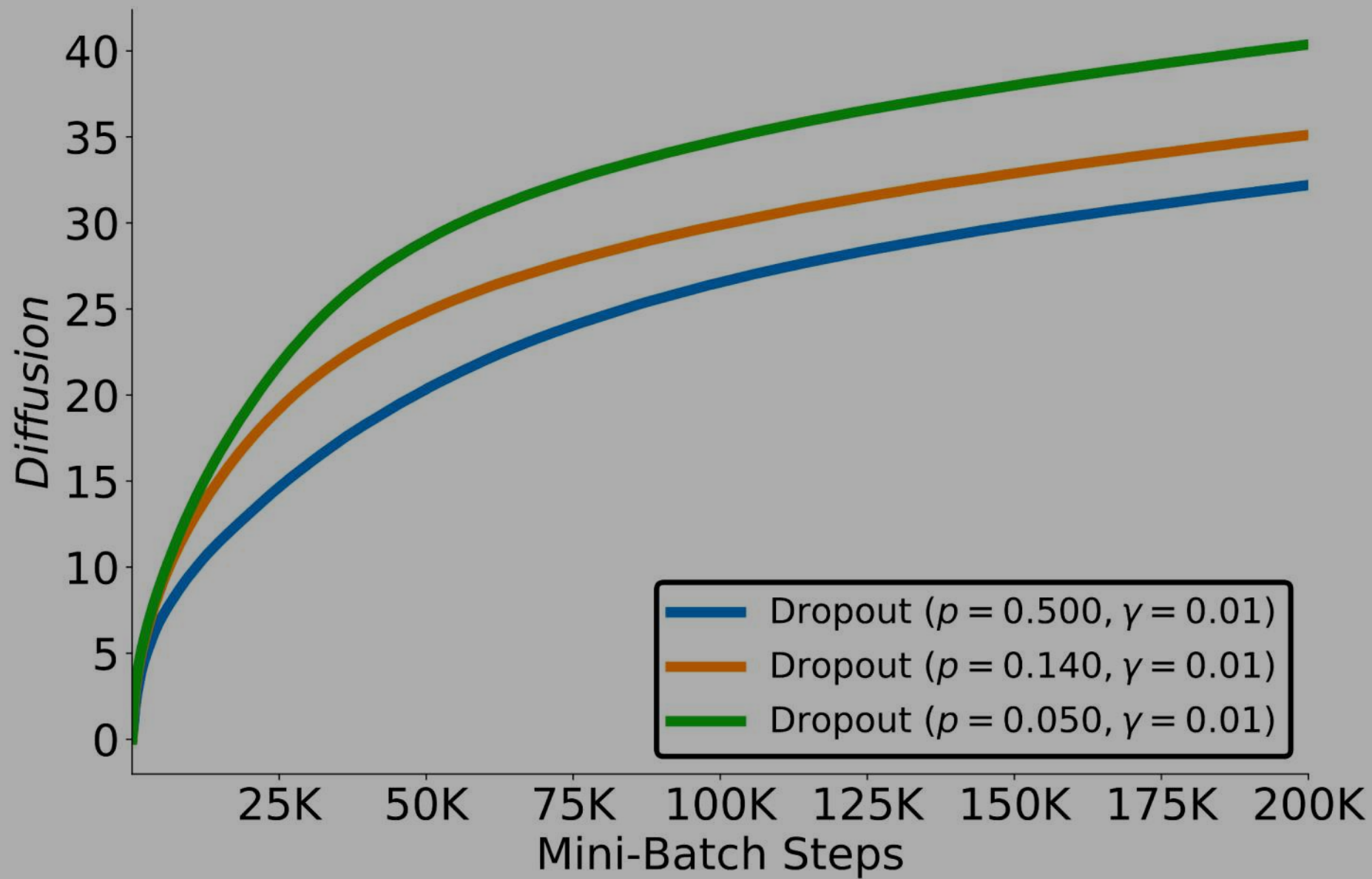
NOISE SPEEDS UP THE DIFFUSION!
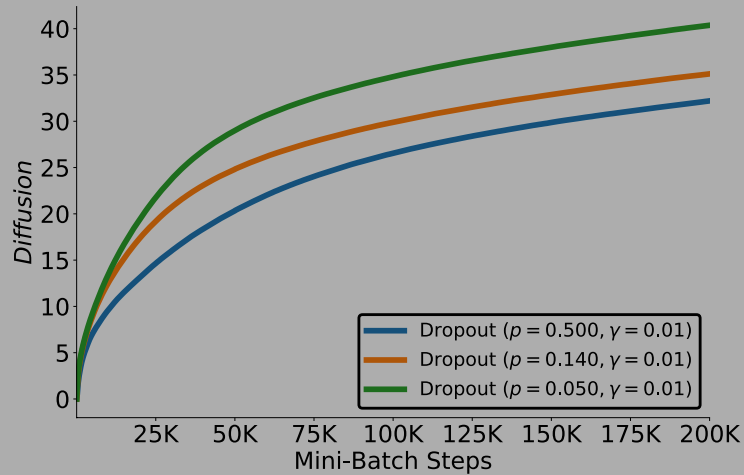
NOISE SPEEDS UP THE DIFFUSION!

# JOINT WEIGHT-HYPERPARAMS SPACE

## GREEDY APPROACH – POPULATION BASED TRAINING



Jaderberg, et.al. arXiv:1711.09846

# JOINT WEIGHT-HYPERPARAMS SPACE

## GREEDY APPROACH – POPULATION BASED TRAINING

GAN population development

FuN population development

4.2   4.5   4.8   5.1   5.4   5.7   6.0   6.3   6.6
Inception Score

1000 2000 3000 4000 5000 6000 7000 8000 9000
Cumulative Expected Reward

Jaderberg, et.al. arXiv:1711.09846

# JOINT WEIGHT-HYPERPARAMS SPACE

## GREEDY APPROACH – POPULATION BASED TRAINING



Jaderberg, et.al. arXiv:1711.09846

# EXPLOITING NOISE-LIKE PROPERTY

ARXIV:1909.04013 (2019)

softserve

# OPTIMIZATION OVER THE PATH



ARXIV:1909.04013 (2019)

softserve

**Algorithm 1** Training with replica exchange

**INPUT:**  Number of replicas $M$, Inverse "temperature" (hyperparameters) $\beta = (\beta_1, \beta_2, \ldots, \beta_M)$
Number of steps for initialization $\Delta N_i$
Number of SGD steps between exchanges $\Delta N_e$
Exchange normalization parameter $C$
Number of steps $T$

**OUTPUT:**  Weight configurations $\mathbf{W} = (\mathbf{W_1}, \mathbf{W_2}, \ldots, \mathbf{W_M})$ of the replicas,
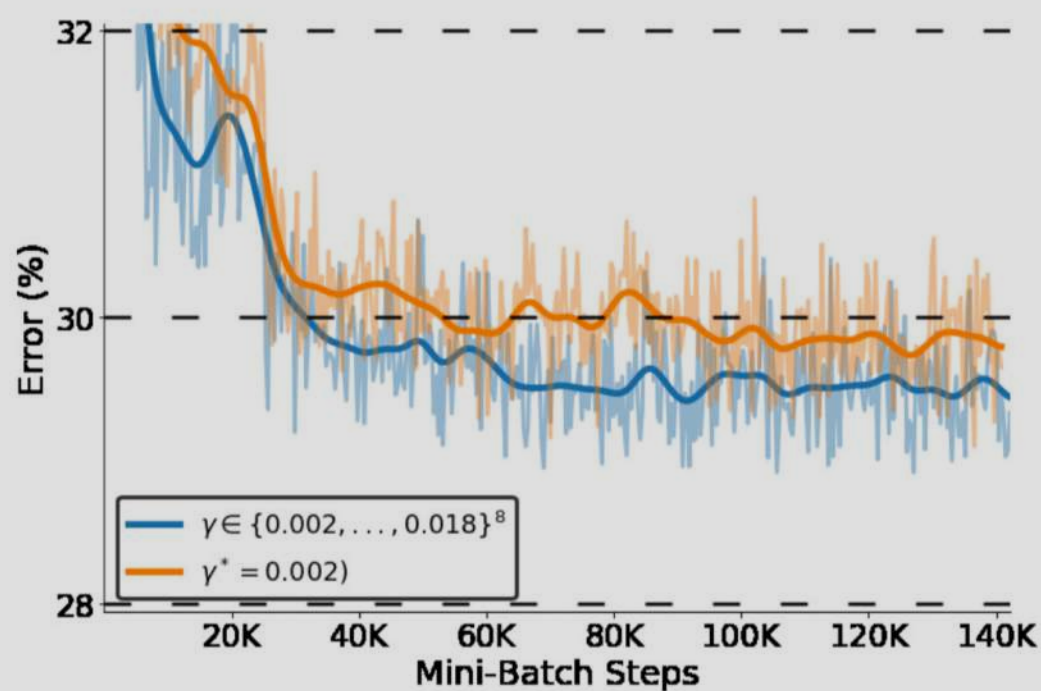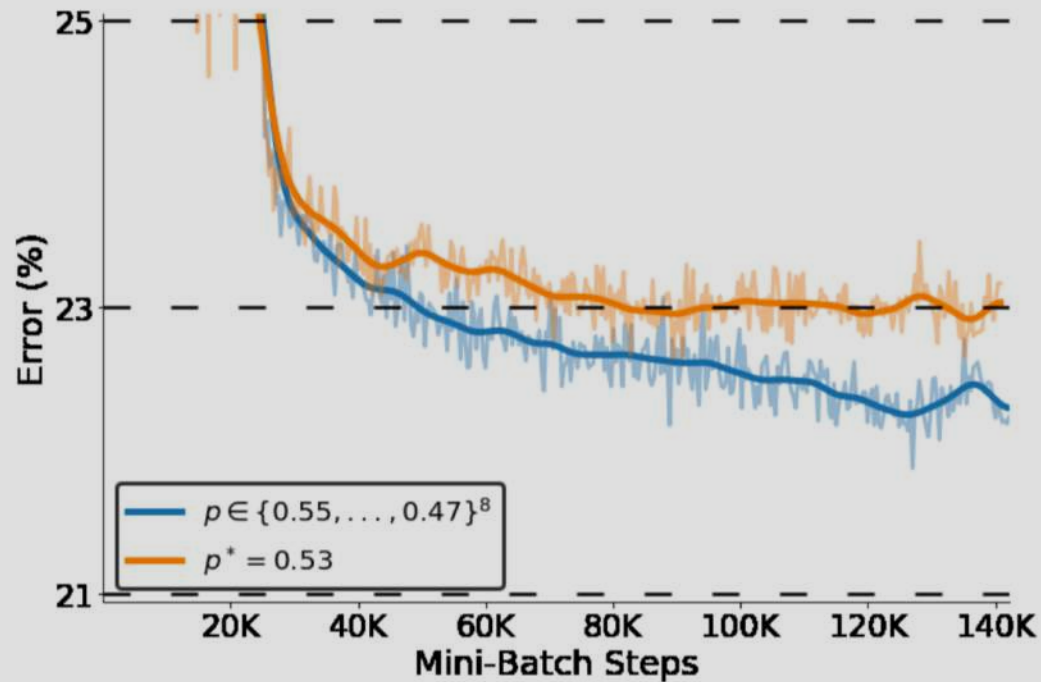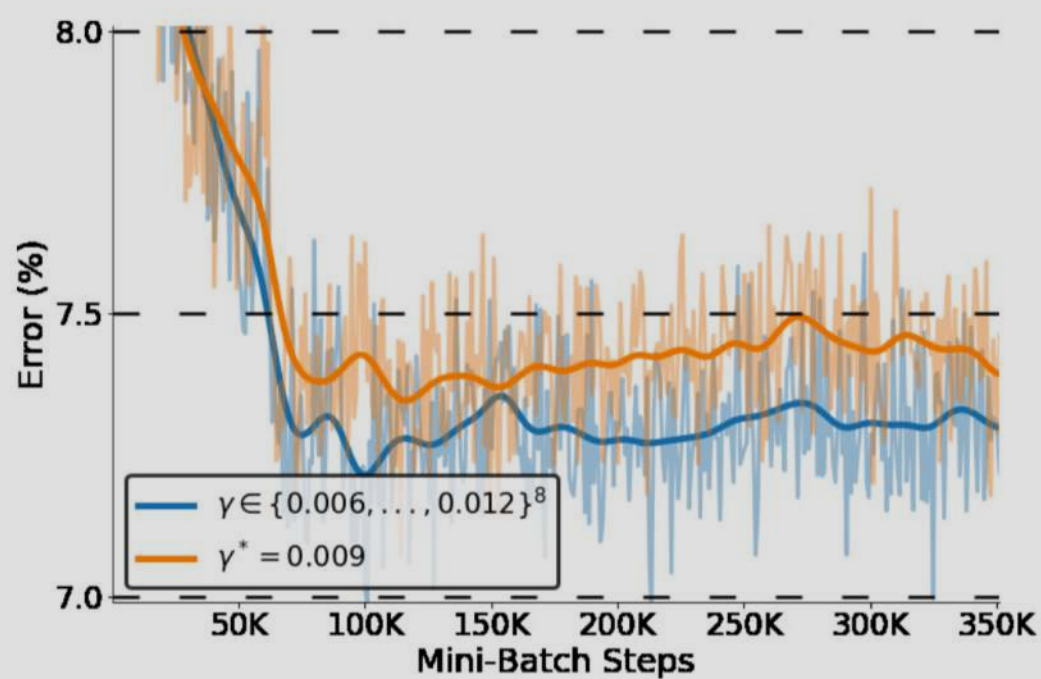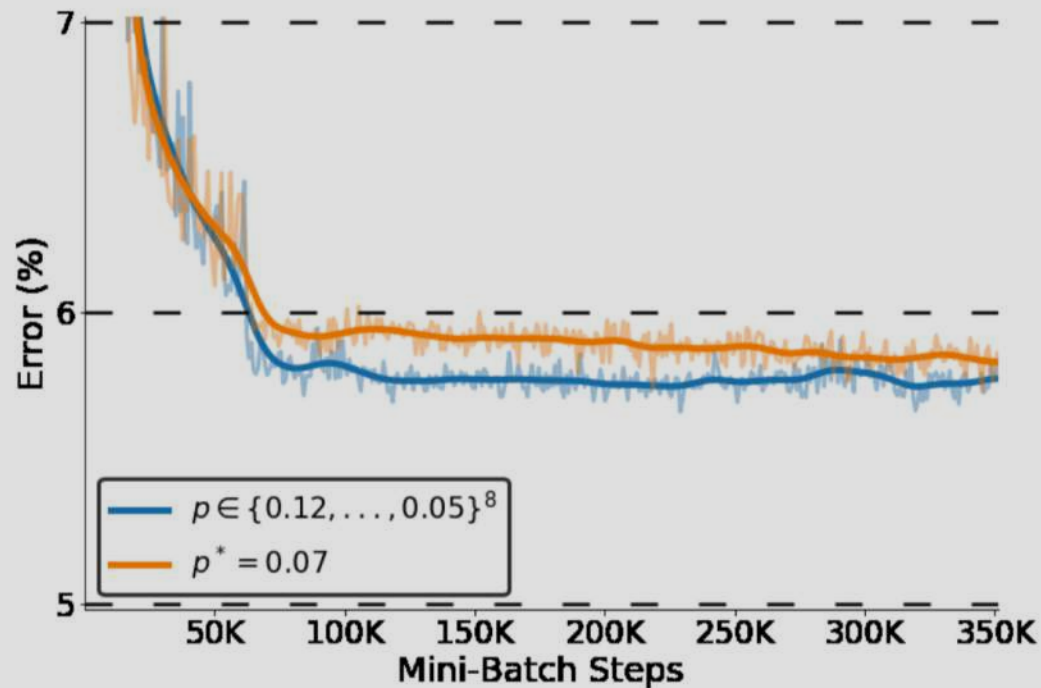  1: **Initialization:** $\forall k \in M$, initialize weights $\mathbf{W}_k$ for each replica and set $t = 0$.
  2: $\forall k \in M$, perform SGD for $\Delta N_i$ steps. Update $t \leftarrow t + 1$ at each step.
  3: **Repeat:**
  4:      $\forall k \in M$, perform SGD for $\Delta N_e$ steps to update $\mathbf{W}_k$. Set $t \leftarrow t + 1$ at each step.
  5:      Let $\mathcal{L}_t = (\mathcal{L}(\mathbf{W}_1^t), \mathcal{L}(\mathbf{W}_2^t), \ldots, \mathcal{L}(\mathbf{W}_k^t))$ be validation losses at time $t$.
  6:      Randomly select a pair $(m, n)$ of replicas with adjacent temperatures.
  7:      **if** $\Delta = C(\beta_m - \beta_n)[\mathcal{L}(\mathbf{W}_m) - \mathcal{L}(\mathbf{W}_n)] \leq 0$ **then**
  8:          swap $\beta_m$ and $\beta_n$
  9:      **else**
 10:          swap $\beta_m$ and $\beta_n$ with probability $\exp(-\Delta)$.
 11:      Update $\alpha$, the acceptance ratio. Finish if $t > T$.

**EMNIST**

Error (%)

25 — 23 — 21

$p \in \{0.55, \ldots, 0.47\}^8$
$p^* = 0.53$

20K  40K  60K  80K  100K  120K  140K
**Mini-Batch Steps**

Error (%)

32 — 30 — 28

$\gamma \in \{0.002, \ldots, 0.018\}^8$
$\gamma^* = 0.002)$

20K  40K  60K  80K  100K  120K  140K
**Mini-Batch Steps**

**CIFAR10**

Error (%)

7 — 6 — 5

$p \in \{0.12, \ldots, 0.05\}^8$
$p^* = 0.07$

50K  100K  150K  200K  250K  300K  350K
**Mini-Batch Steps**

Error (%)

8.0 — 7.5 — 7.0

$\gamma \in \{0.006, \ldots, 0.012\}^8$
$\gamma^* = 0.009$

50K  100K  150K  200K  250K  300K  350K
**Mini-Batch Steps**

# RESNET/CIFAR10

ARXIV:1909.04013 (2019)

WITH

VLAD PUSHKARIOV, TECHNION
YONATHAN EFRONI, TECHNION
MACIEJ KOCH-JANUSZ, ETH ZURICH

- TRAINING IN A JOINT WEIGHTS-HYPERPARAMETER SPACE
- EFFICIENT SAMPLING OF THE HYPERPARAMETERS
- NON-STATIONARY SCHEDULING PROTOCOL
- FASTER AND BETTER TRAINED MODELS

# THANKS FOR ATTENTION

Mykola Maksymenko

softserve