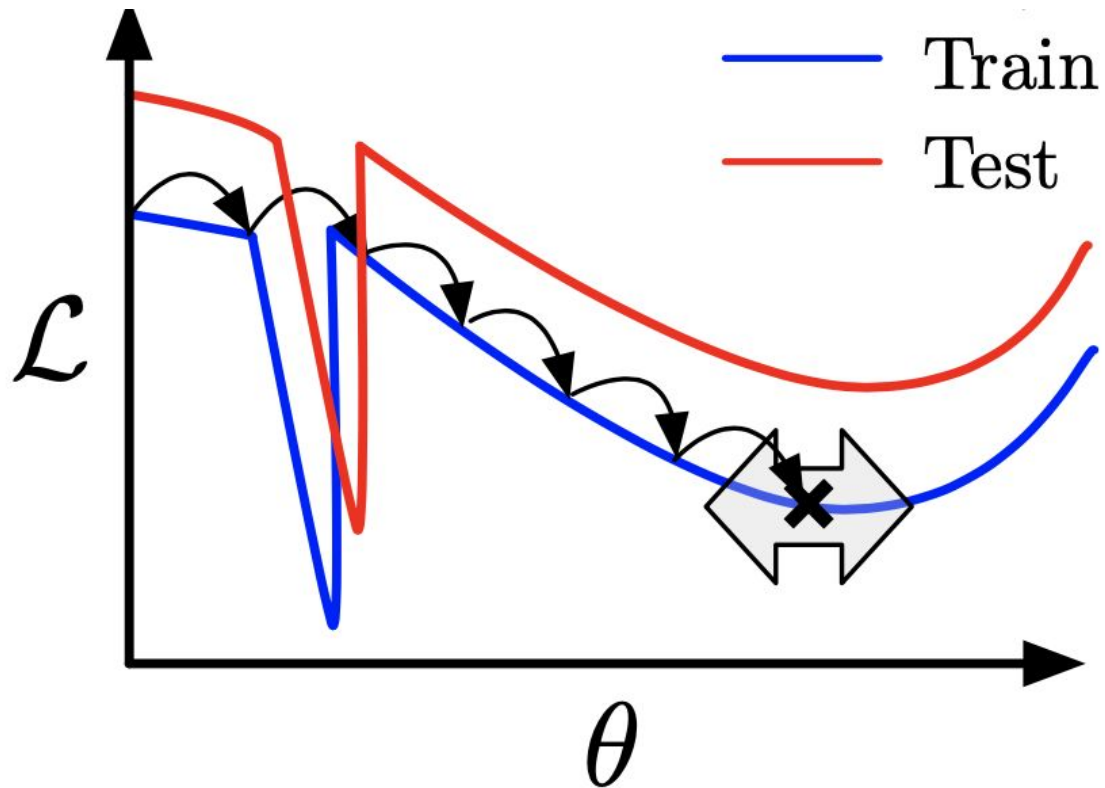


# Широкие и узкие оптимумы у нейросетей

Писцов Георгий, БПМИ202

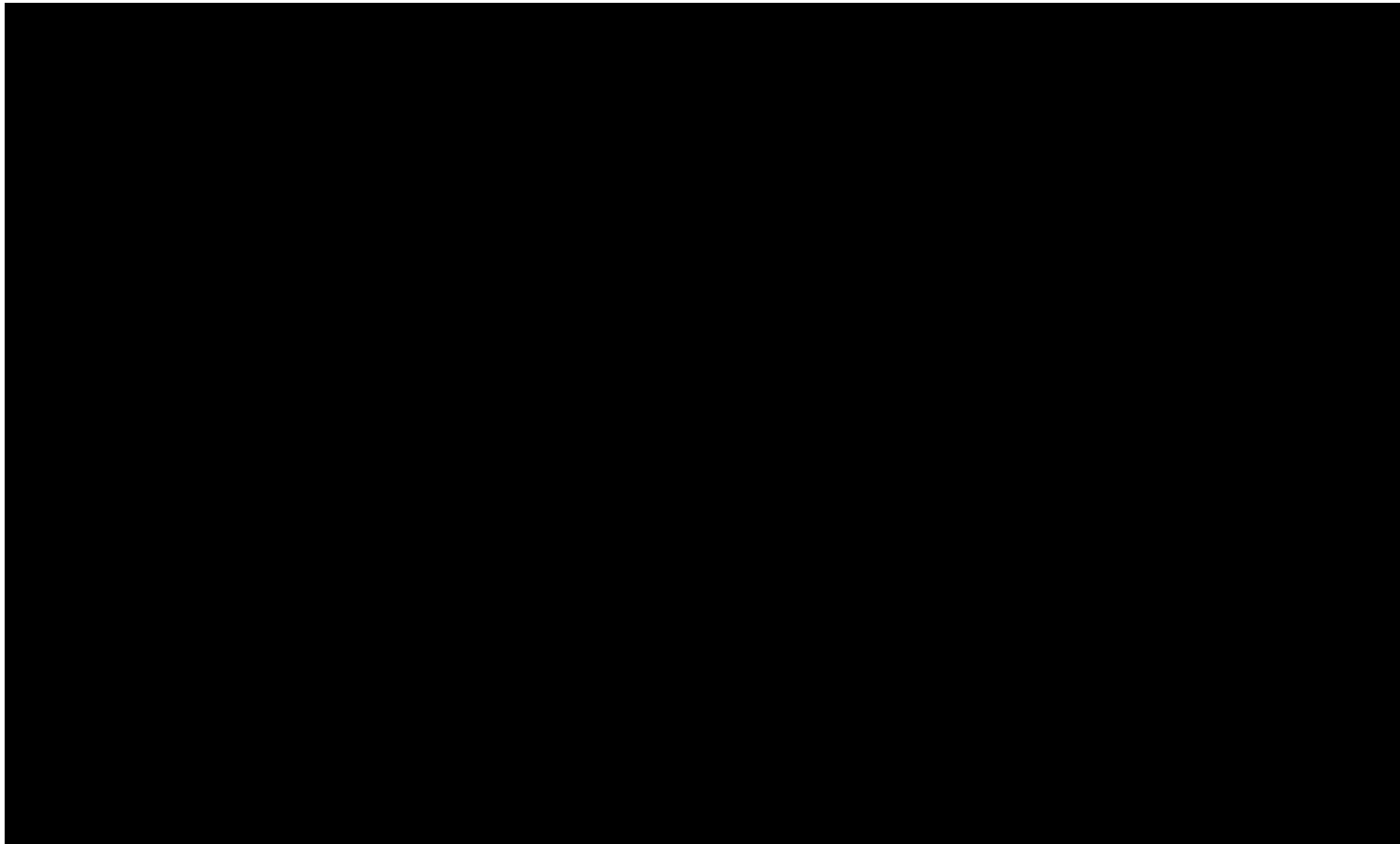
# Что такое широкие и узкие оптимумы?



Идея: Сходимость к широким оптимумам ведёт к лучшей обобщенности (generalization) модели, так как при небольшом возмущении значение остается приблизительно оптимальным

В данной работе мы рассмотрим два алгоритма, которые позволяют нам искать широкие оптимумы: SWA и SAM

# Наглядная визуализация узких и широких оптимумов



# Анализ SGD

Мы рассматриваем циклический или константный шаг(вырожденный), заданный формулой:

$$\alpha(i) = (1 - t(i))\alpha_1 + t(i)\alpha_2,$$
$$t(i) = \frac{1}{c} (\text{mod}(i - 1, c) + 1).$$

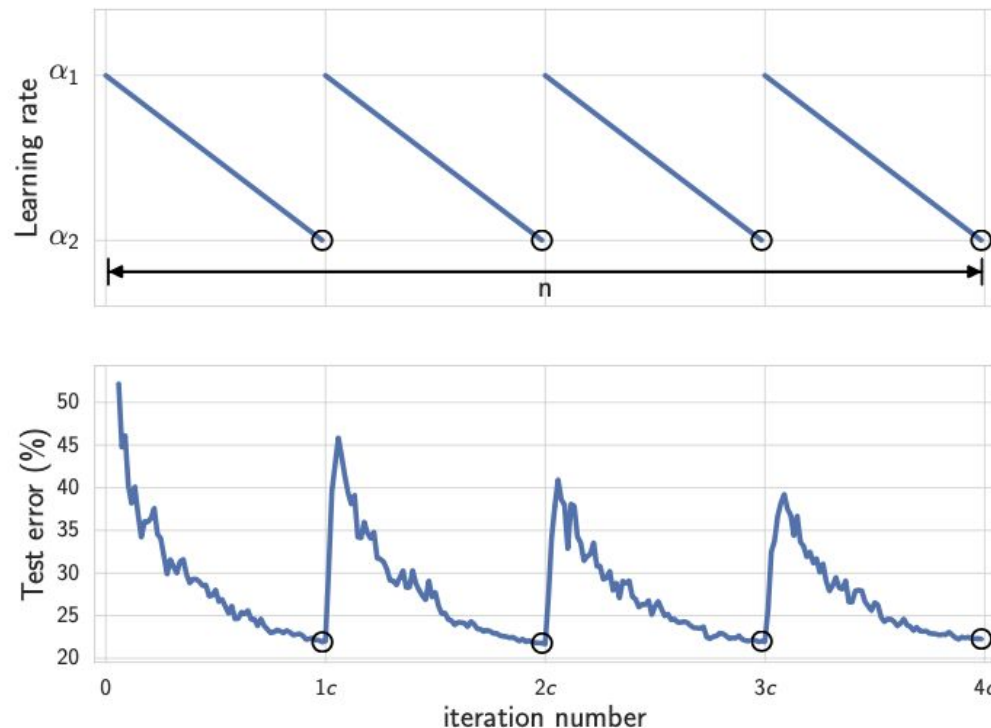


Figure 2: **Top:** cyclical learning rate as a function of iteration. **Bottom:** test error as a function of iteration for cyclical learning rate schedule with Preactivation-ResNet-164 on CIFAR-100. Circles indicate iterations corresponding to the minimum learning rates.

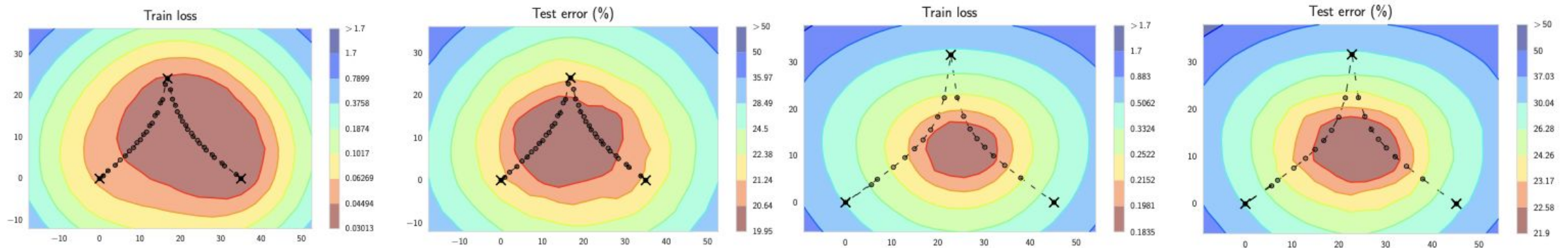


Figure 3: The  $L_2$ -regularized cross-entropy train loss and test error surfaces of a Preactivation ResNet-164 on CIFAR-100 in the plane containing the first, middle and last points (indicated by black crosses) in the trajectories with **(left two)** cyclical and **(right two)** constant learning rate schedules.

Замечаем, что:

1. Оба случая исследуют точки близко к периферии набора высокоэффективных сетей
2. Есть тот самый сдвиг между поверхностью потерь на тесте и на обучении
3. Усреднение наталкивает на мысль о возможном успехе

# Алгоритм SWA

---

**Algorithm 1** Stochastic Weight Averaging

---

**Require:**

weights  $\hat{w}$ , LR bounds  $\alpha_1, \alpha_2$ ,  
cycle length  $c$  (for constant learning rate  $c = 1$ ), number of iterations  $n$

**Ensure:**  $w_{\text{SWA}}$ 

$w \leftarrow \hat{w}$  {Initialize weights with  $\hat{w}$ }

$w_{\text{SWA}} \leftarrow w$

**for**  $i \leftarrow 1, 2, \dots, n$  **do**

$\alpha \leftarrow \alpha(i)$  {Calculate LR for the iteration}

$w \leftarrow w - \alpha \nabla \mathcal{L}_i(w)$  {Stochastic gradient update}

**if**  $\text{mod}(i, c) = 0$  **then**

$n_{\text{models}} \leftarrow i/c$  {Number of models}

$w_{\text{SWA}} \leftarrow \frac{w_{\text{SWA}} \cdot n_{\text{models}} + w}{n_{\text{models}} + 1}$  {Update average}

**end if**

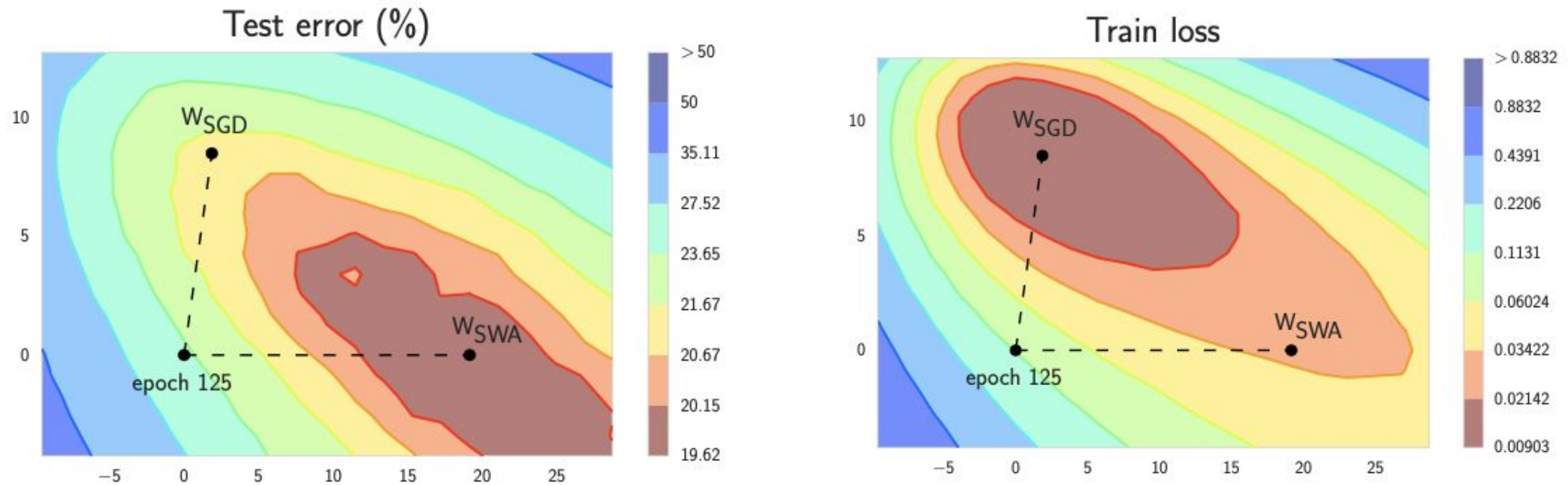
**end for**

{Compute BatchNorm statistics for  $w_{\text{SWA}}$  weights}

---

Чуть позже мы увидим более современную версию алгоритма, которая вносит небольшую модификацию

# Сравнение SWA и SGD



Иллюстрации SWA и SGD с преактивацией ResNet-164 на CIFAR-100. Слева и справа: поверхности ошибок теста и потерь обучения, показывающие веса, предложенные SGD (при сходимости) и SWA, начиная с той же инициализации SGD после 125 эпох обучения.



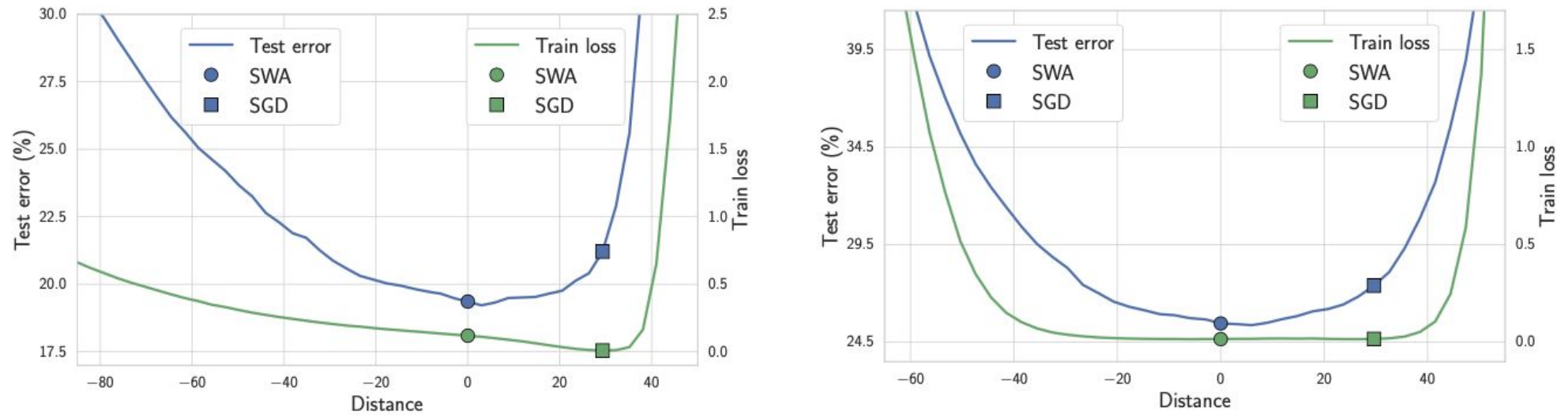


Figure 5:  $L_2$ -regularized cross-entropy train loss and test error as a function of a point on the line connecting SWA and SGD solutions on CIFAR-100. **Left:** Preactivation ResNet-164. **Right:** VGG-16.

Рассматриваем линейную интерполяцию::  $w(t) = t \cdot w_{\text{SGD}} + (1 - t) \cdot w_{\text{SWA}}$

Замечаем, что:

1. Графики ошибок действительно смещены
2. SWA находит более широкий оптимум
3. SGD находит оптимум, в одном из направлений которого ошибка резко растёт



# Связь SWA и FGE: теория

- В первом случае усреднение весов
- Во втором случае усреднение прогнозов

$$\begin{aligned}\bar{f} &= \frac{1}{n} \sum_{i=1}^n f(w_i). & f(w_j) &= f(w_{\text{SWA}}) + \langle \nabla f(w_{\text{SWA}}), \Delta_j \rangle + O(\|\Delta_j\|^2), \\ \Delta_i &= w_i - w_{\text{SWA}} & \bar{f} - f(w_{\text{SWA}}) &= \frac{1}{n} \sum_{i=1}^n (\langle \nabla f(w_{\text{SWA}}), \Delta_i \rangle + O(\|\Delta_i\|^2)) \\ \Delta &= \max_{i=1}^n \|\Delta_i\| & &= \left\langle \nabla f(w_{\text{SWA}}), \frac{1}{n} \sum_{i=1}^n \Delta_i \right\rangle + O(\Delta^2) = O(\Delta^2),\end{aligned}$$

# Связь SWA и FGE: практика

Мы запускаем FGE на 20 эпох и сравнили предсказания различных моделей на тестовом наборе данных с преактивацией ResNet-164 на CIFAR-100. Затем мы усредняем веса предложений и вычисляем вероятности классов на тестовом наборе данных. Для FGE и SWA доля идентично помеченных объектов составляет 95,26%.

# Сравнение SGD, FGE и SWA

Table 1: Accuracies (%) of SWA, SGD and FGE methods on CIFAR-100 and CIFAR-10 datasets for different training budgets. Accuracies for the FGE ensemble are from [Garipov et al. \[2018\]](#).

DNN (Budget)	SGD	FGE (1 Budget)	SWA		
			1 Budget	1.25 Budgets	1.5 Budgets
CIFAR-100					
VGG-16 (200)	72.55 ± 0.10	74.26	73.91 ± 0.12	74.17 ± 0.15	74.27 ± 0.25
ResNet-164 (150)	78.49 ± 0.36	79.84	79.77 ± 0.17	80.18 ± 0.23	80.35 ± 0.16
WRN-28-10 (200)	80.82 ± 0.23	82.27	81.46 ± 0.23	81.91 ± 0.27	82.15 ± 0.27
PyramidNet-272 (300)	83.41 ± 0.21	—	—	83.93 ± 0.18	84.16 ± 0.15
CIFAR-10					
VGG-16 (200)	93.25 ± 0.16	93.52	93.59 ± 0.16	93.70 ± 0.22	93.64 ± 0.18
ResNet-164 (150)	95.28 ± 0.10	95.45	95.56 ± 0.11	95.77 ± 0.04	95.83 ± 0.03
WRN-28-10 (200)	96.18 ± 0.11	96.36	96.45 ± 0.11	96.64 ± 0.08	96.79 ± 0.05
ShakeShake-2x64d (1800)	96.93 ± 0.10	—	—	97.16 ± 0.10	97.12 ± 0.06

# Расписание шага SWA

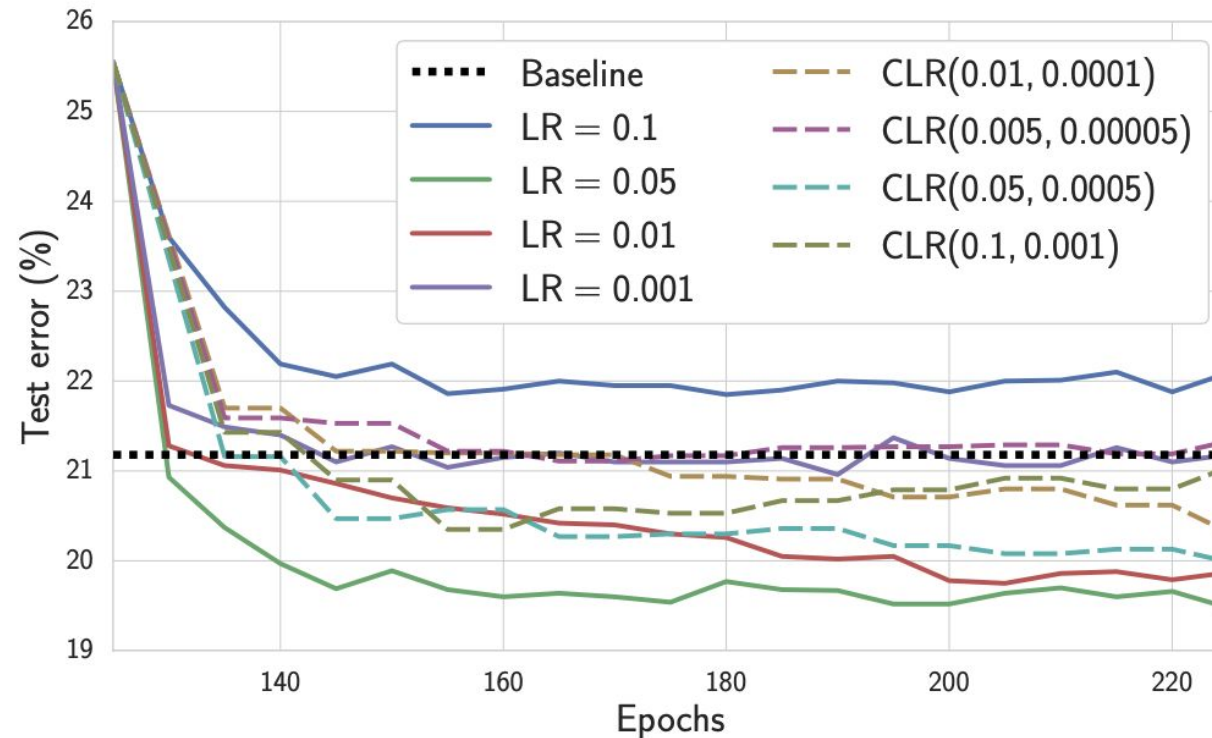


Figure 6: Test error as a function of training epoch for SWA with different learning rate schedules with a Preactivation ResNet-164 on CIFAR-100.

# Другие подходы: SAM

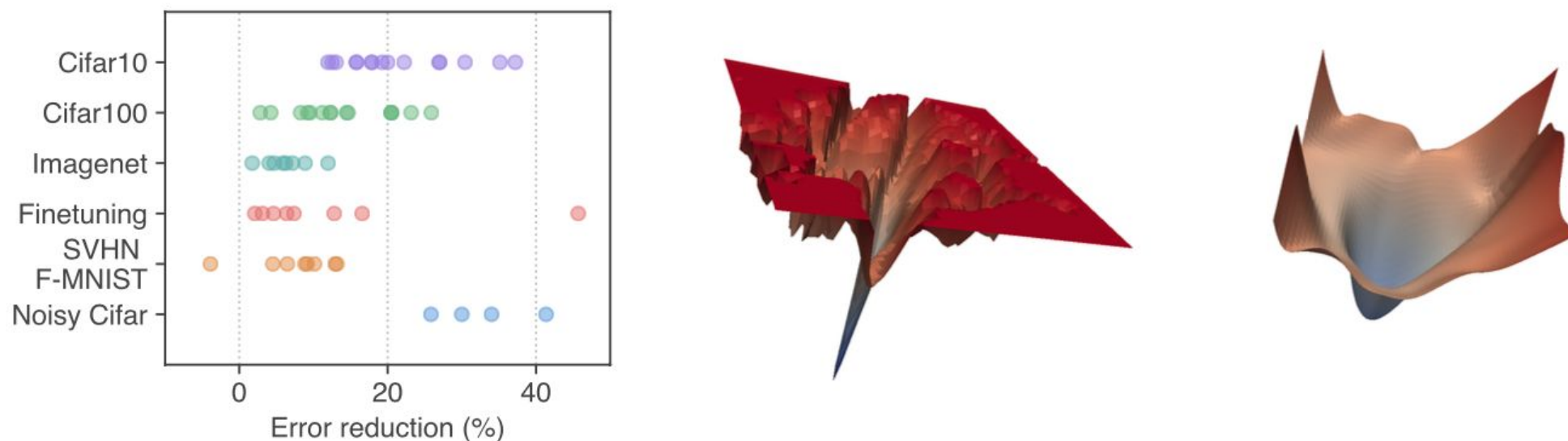


Figure 1: (left) Error rate reduction obtained by switching to SAM. Each point is a different dataset / model / data augmentation. (middle) A sharp minimum to which a ResNet trained with SGD converged. (right) A wide minimum to which the same ResNet trained with SAM converged.

# Описание алгоритма SAM

Идея состоит в том, что вместо того, чтобы искать значения параметров  $w$ , которые просто имеют низкое значение потерь при обучении, мы ищем значения параметров, все окрестности которых имеют равномерно низкое значение потерь при обучении. Для этого мы будем решать минимаксную задачу

$$\min_{\theta} \max_{\|\epsilon\|_2 \leq \rho} \mathcal{L}(\theta + \epsilon),$$

$$\hat{\epsilon}(w) = \rho \operatorname{sign}(\nabla_w L_{\mathcal{S}}(w)) |\nabla_w L_{\mathcal{S}}(w)|^{q-1} / \left( \|\nabla_w L_{\mathcal{S}}(w)\|_q^q \right)^{1/p}$$

$$\nabla_w L_{\mathcal{S}}^{SAM}(w) \approx \nabla_w L_{\mathcal{S}}(w)|_{w+\hat{\epsilon}(w)}.$$



# SAM реализация

**Input:** Training set  $\mathcal{S} \triangleq \cup_{i=1}^n \{(\mathbf{x}_i, \mathbf{y}_i)\}$ , Loss function  $l : \mathcal{W} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ , Batch size  $b$ , Step size  $\eta > 0$ , Neighborhood size  $\rho > 0$ .

**Output:** Model trained with SAM

Initialize weights  $\mathbf{w}_0, t = 0$ ;

**while** *not converged* **do**

    Sample batch  $\mathcal{B} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_b, \mathbf{y}_b)\}$ ;

    Compute gradient  $\nabla_{\mathbf{w}} L_{\mathcal{B}}(\mathbf{w})$  of the batch's training loss;

    Compute  $\hat{\epsilon}(\mathbf{w})$  per equation 2;

    Compute gradient approximation for the SAM objective

        (equation 3):  $\mathbf{g} = \nabla_{\mathbf{w}} L_{\mathcal{B}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})}$ ;

    Update weights:  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}$ ;

$t = t + 1$ ;

**end**

**return**  $\mathbf{w}_t$

**Algorithm 1:** SAM algorithm

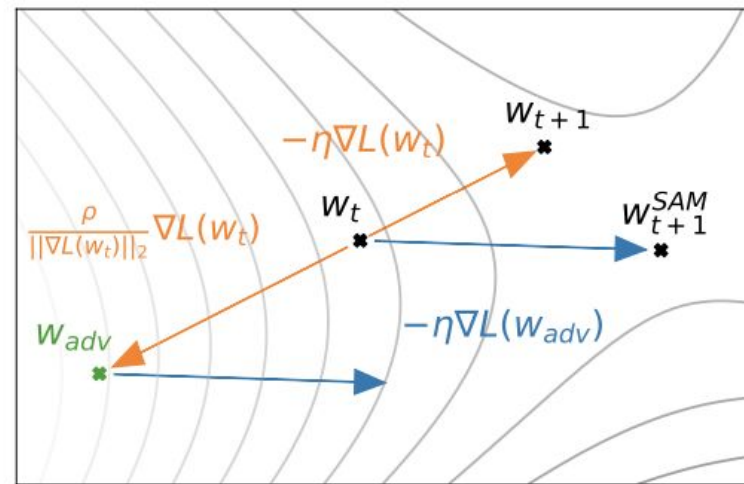


Figure 2: Schematic of the SAM parameter update.

# Сравнение SAM с SGD на примерах

Model	Augmentation	CIFAR-10		CIFAR-100	
		SAM	SGD	SAM	SGD
WRN-28-10 (200 epochs)	Basic	<b>2.7</b> $\pm 0.1$	3.5 $\pm 0.1$	<b>16.5</b> $\pm 0.2$	18.8 $\pm 0.2$
WRN-28-10 (200 epochs)	Cutout	<b>2.3</b> $\pm 0.1$	2.6 $\pm 0.1$	<b>14.9</b> $\pm 0.2$	16.9 $\pm 0.1$
WRN-28-10 (200 epochs)	AA	<b>2.1</b> $\pm <0.1$	2.3 $\pm 0.1$	<b>13.6</b> $\pm 0.2$	15.8 $\pm 0.2$
WRN-28-10 (1800 epochs)	Basic	<b>2.4</b> $\pm 0.1$	3.5 $\pm 0.1$	<b>16.3</b> $\pm 0.2$	19.1 $\pm 0.1$
WRN-28-10 (1800 epochs)	Cutout	<b>2.1</b> $\pm 0.1$	2.7 $\pm 0.1$	<b>14.0</b> $\pm 0.1$	17.4 $\pm 0.1$
WRN-28-10 (1800 epochs)	AA	<b>1.6</b> $\pm 0.1$	2.2 $\pm <0.1$	<b>12.8</b> $\pm 0.2$	16.1 $\pm 0.2$
Shake-Shake (26 2x96d)	Basic	<b>2.3</b> $\pm <0.1$	2.7 $\pm 0.1$	<b>15.1</b> $\pm 0.1$	17.0 $\pm 0.1$
Shake-Shake (26 2x96d)	Cutout	<b>2.0</b> $\pm <0.1$	2.3 $\pm 0.1$	<b>14.2</b> $\pm 0.2$	15.7 $\pm 0.2$
Shake-Shake (26 2x96d)	AA	<b>1.6</b> $\pm <0.1$	1.9 $\pm 0.1$	<b>12.8</b> $\pm 0.1$	14.1 $\pm 0.2$
PyramidNet	Basic	<b>2.7</b> $\pm 0.1$	4.0 $\pm 0.1$	<b>14.6</b> $\pm 0.4$	19.7 $\pm 0.3$
PyramidNet	Cutout	<b>1.9</b> $\pm 0.1$	2.5 $\pm 0.1$	<b>12.6</b> $\pm 0.2$	16.4 $\pm 0.1$
PyramidNet	AA	<b>1.6</b> $\pm 0.1$	1.9 $\pm 0.1$	<b>11.6</b> $\pm 0.1$	14.6 $\pm 0.1$
PyramidNet+ShakeDrop	Basic	<b>2.1</b> $\pm 0.1$	2.5 $\pm 0.1$	<b>13.3</b> $\pm 0.2$	14.5 $\pm 0.1$
PyramidNet+ShakeDrop	Cutout	<b>1.6</b> $\pm <0.1$	1.9 $\pm 0.1$	<b>11.3</b> $\pm 0.1$	11.8 $\pm 0.2$
PyramidNet+ShakeDrop	AA	<b>1.4</b> $\pm <0.1$	1.6 $\pm <0.1$	<b>10.3</b> $\pm 0.1$	10.6 $\pm 0.1$

Table 1: Results for SAM on state-of-the-art models on CIFAR- $\{10, 100\}$  (WRN = WideResNet; AA = AutoAugment; SGD is the standard non-SAM procedure used to train these models).

# УСТОЙЧИВОСТЬ К ШУМУ

Method	Noise rate (%)			
	20	40	60	80
<a href="#">Sanchez et al. (2019)</a>	94.0	92.8	90.3	74.1
<a href="#">Zhang &amp; Sabuncu (2018)</a>	89.7	87.6	82.7	67.9
<a href="#">Lee et al. (2019)</a>	87.1	81.8	75.4	-
<a href="#">Chen et al. (2019)</a>	89.7	-	-	52.3
<a href="#">Huang et al. (2019)</a>	92.6	90.3	43.4	-
MentorNet (2017)	92.0	91.2	74.2	60.0
Mixup (2017)	94.0	91.5	86.8	76.9
MentorMix (2019)	<b>95.6</b>	<b>94.2</b>	91.3	<b>81.0</b>
SGD	84.8	68.8	48.2	26.2
Mixup	93.0	90.0	83.8	70.2
Bootstrap + Mixup	93.3	92.0	87.6	72.0
SAM	95.1	93.4	90.5	77.9
Bootstrap + SAM	95.4	<b>94.2</b>	<b>91.8</b>	79.9

Table 4: Test accuracy on the clean test set for models trained on CIFAR-10 with noisy labels. Lower block is our implementation, upper block gives scores from the literature, per [Jiang et al. \(2019\)](#).

Тот факт, что SAM ищет параметры модели, устойчивые к возмущениям, позволяет предположить, что SAM может быть устойчив к шуму в обучающем множестве (который может возмутить ландшафт потерь при обучении).

# НОВЫЕ ПОДХОДЫ

---

**Algorithm 1** Stochastic Weight Averaging [48]

**Input:** Loss function  $\mathcal{L}$ , training budget in number of iterations  $b$ , training dataset  $\mathcal{D} := \cup_{i=1}^n \{\mathbf{x}_i\}$ , mini-batch size  $|\mathcal{B}|$ , averaging start epoch  $E$ , averaging frequency  $\nu$ , (scheduled) learning rate  $\eta$ , initial weights  $\theta_0$ .

```
for  $k \leftarrow 1, \dots, b$  do
  Sample a mini-batch  $\mathcal{B}$  from  $\mathcal{D}$ 
  Compute gradient  $\mathbf{g} \leftarrow \nabla \mathcal{L}(\theta_t)$ 
  Update parameters  $\theta_{t+1} \leftarrow \theta_t - \eta \mathbf{g}$ 
  if  $k \geq E$  and  $\text{mod}(k, \nu) = 0$  then
     $\theta_{t+1}^{\text{SWA}} = (\theta_t^{\text{SWA}} \cdot l + \theta_{t+1}^{\text{SWA}}) / (l + 1)$ 
  end if
end for
return  $\theta^{\text{SWA}}$ 
```

---

---

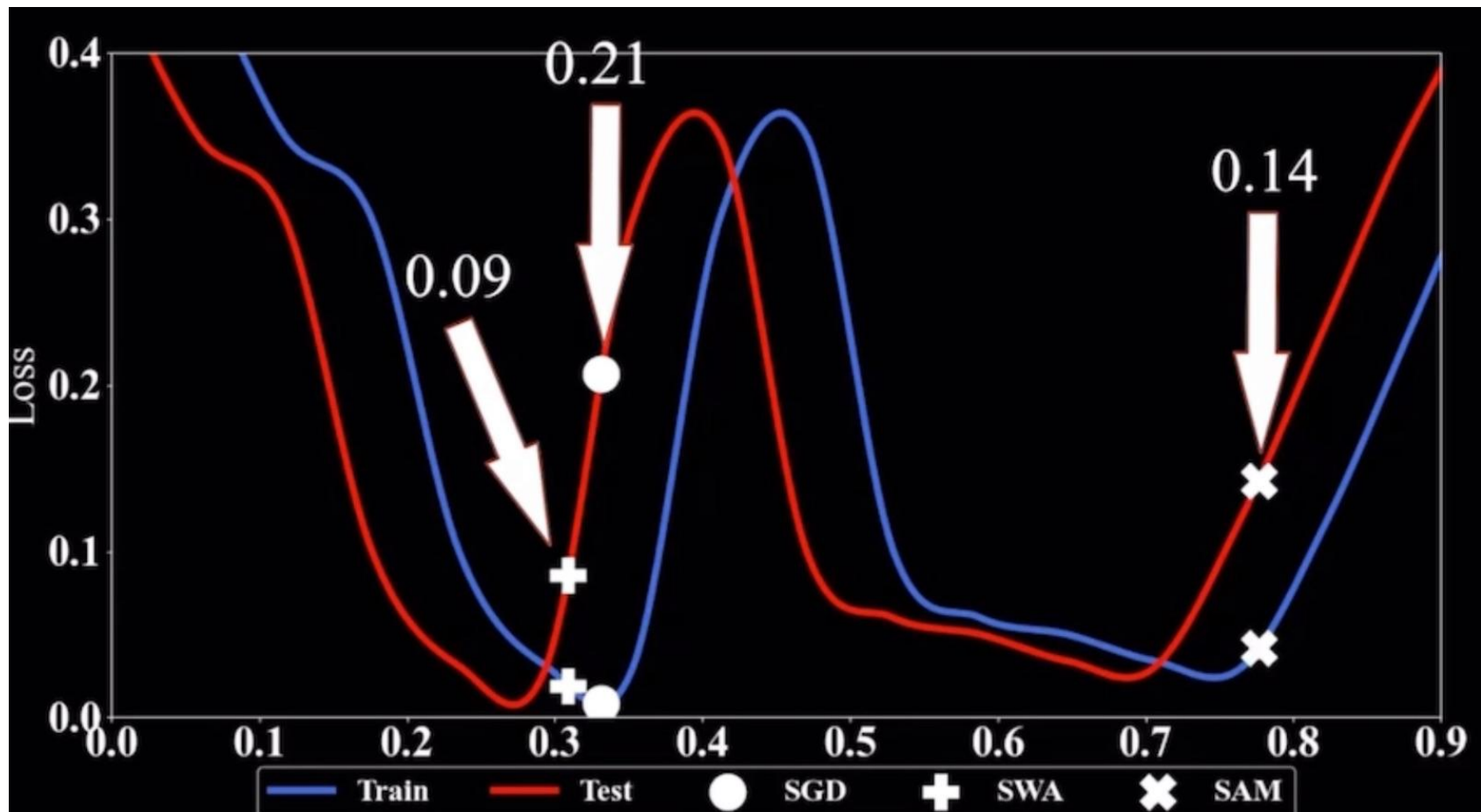
**Algorithm 2** Sharpness-Aware Minimization [22]

**Input:** Loss function  $\mathcal{L}$ , training budget in number of iterations  $b$ , training dataset  $\mathcal{D} := \cup_{i=1}^n \{\mathbf{x}_i\}$ , mini-batch size  $|\mathcal{B}|$ , neighborhood radius  $\rho$ , (scheduled) learning rate  $\eta$ , initial weights  $\theta_0$ .

```
for  $k \leftarrow 1, \dots, b$  do
  Sample a mini-batch  $\mathcal{B}$  from  $\mathcal{D}$ 
  Compute worst-case perturbation
     $\hat{\epsilon} \leftarrow \rho \frac{\nabla \mathcal{L}(\theta)}{\|\nabla \mathcal{L}(\theta)\|_2}$ 
  Compute gradient  $\mathbf{g} \leftarrow \nabla \mathcal{L}(\theta_t^{\text{SAM}} + \hat{\epsilon})$ 
  Update parameters  $\theta_{t+1}^{\text{SAM}} \leftarrow \theta_t^{\text{SAM}} - \eta \mathbf{g}$ 
end for
return  $\theta^{\text{SAM}}$ 
```

---

# Чем же отличаются оптимумы SAM и SWA





# Чем же отличаются оптимумы SAM и SWA

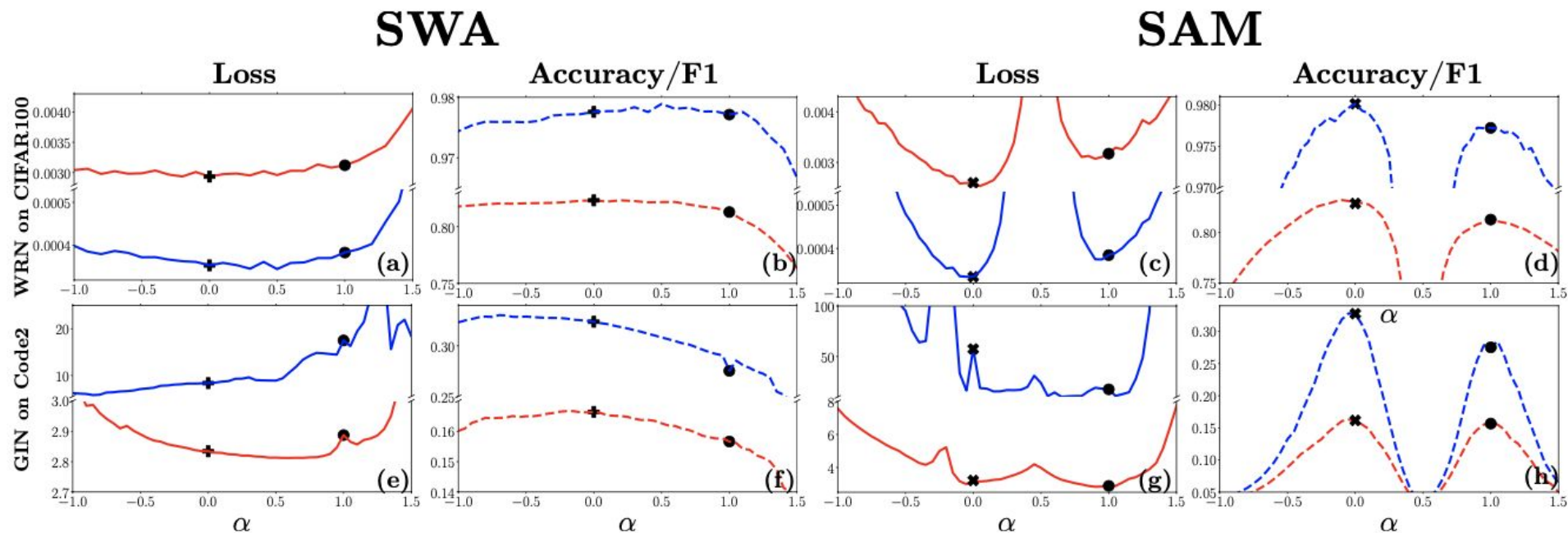


Figure 2: Training (blue) and test (red) losses (—) and accuracies (-----) of linear interpolations  $\theta(\alpha) = (1 - \alpha)\theta + \alpha\theta'$  (for  $\alpha \in [-1, 1.5]$ ) between SWA (+) and SAM (x) solutions ( $\alpha = 0.0$ ) and non-flat baseline solutions ( $\bullet$ ,  $\alpha = 1.0$ ).



# Сравниваем все на примерах

Table 2: CV test results: Supervised Classification (SC), and Self-Supervised Learning (SSL) tasks.

Task	Model	Baseline	SWA	SAM	SWA+SAM
SC: CIFAR10	WRN-28-10	96.78 $\pm$ 0.03	-0.05 $\pm$ 0.04	+ 0.34 $\pm$ 0.09	+ 0.25 $\pm$ 0.05
	PN-272	96.73 $\pm$ 0.14	+ 0.22 $\pm$ 0.14	+ 0.42 $\pm$ 0.06	+ 0.41 $\pm$ 0.02
	ViT-B-16	98.95 $\pm$ 0.02	-0.04 $\pm$ 0.04	+ 0.07 $\pm$ 0.01	+ 0.10 $\pm$ 0.01
	Mixer-B-16	96.65 $\pm$ 0.03	+ 0.02 $\pm$ 0.03	+ 0.19 $\pm$ 0.05	+ 0.22 $\pm$ 0.06
SC: CIFAR100	WRN-28-10	80.93 $\pm$ 0.19	+ 1.62 $\pm$ 0.06	+ 1.82 $\pm$ 0.14	+ 2.24 $\pm$ 0.14
	PN-272	80.86 $\pm$ 0.12	+ 1.88 $\pm$ 0.04	+ 2.33 $\pm$ 0.08	+ 2.60 $\pm$ 0.09
	ViT-B-16	92.77 $\pm$ 0.07	-0.12 $\pm$ 0.05	+ 0.19 $\pm$ 0.09	+ 0.13 $\pm$ 0.07
	Mixer-B-16	83.77 $\pm$ 0.08	+ 0.45 $\pm$ 0.06	+ 0.52 $\pm$ 0.15	+ 0.97 $\pm$ 0.12
SSL: CIFAR10	MoCo	89.25 $\pm$ 0.07	-0.03 $\pm$ 0.10	-0.25 $\pm$ 0.06	-0.17 $\pm$ 0.10
	SimCLR	88.66 $\pm$ 0.08	-0.05 $\pm$ 0.06	+ 0.05 $\pm$ 0.04	-0.13 $\pm$ 0.06
	SimSiam	89.86 $\pm$ 0.22	+ 0.12 $\pm$ 0.26	+ 0.07 $\pm$ 0.10	+ 0.11 $\pm$ 0.10
	BarlowTwins	86.34 $\pm$ 0.24	-0.09 $\pm$ 0.19	+ 0.09 $\pm$ 0.15	+ 0.14 $\pm$ 0.05
	BYOL	90.32 $\pm$ 0.14	+ 0.70 $\pm$ 0.05	+ 0.14 $\pm$ 0.03	+ 0.21 $\pm$ 0.07
	SwaV	87.28 $\pm$ 0.05	+ 0.09 $\pm$ 0.06	+ 0.07 $\pm$ 0.12	+ 0.02 $\pm$ 0.06
SSL: ImageNette	MoCo	81.74 $\pm$ 0.18	+ 0.97 $\pm$ 0.10	+ 0.91 $\pm$ 0.32	+ 1.40 $\pm$ 0.10
	SimCLR	83.28 $\pm$ 0.22	+ 0.95 $\pm$ 0.25	+ 0.18 $\pm$ 0.24	+ 1.07 $\pm$ 0.13
	SimSiam	81.77 $\pm$ 0.14	+ 0.20 $\pm$ 0.37	+ 0.33 $\pm$ 0.28	+ 0.18 $\pm$ 0.26
	BarlowTwins	77.49 $\pm$ 0.36	+ 0.20 $\pm$ 0.16	+ 0.47 $\pm$ 0.27	+ 0.66 $\pm$ 0.57
	BYOL	84.16 $\pm$ 0.14	+ 0.76 $\pm$ 0.08	+ 0.15 $\pm$ 0.25	+ 0.31 $\pm$ 0.19
	SwaV	88.16 $\pm$ 0.31	+ 1.04 $\pm$ 0.27	+ 0.03 $\pm$ 0.10	+ 1.03 $\pm$ 0.09

# Сравниваем все на примерах

Figure 4: (a) NLP test results: Open-Domain Question Answering and Natural Language Understanding (GLUE) including paraphrase, sentiment analysis, and textual entailment. (b) GRL test results: Node Property Prediction (NPP), Graph Property Prediction (GPP), Link Property Prediction (LPP).

Task	Model	Baseline	SWA	SAM	SWA+SAM
NQ	FiD	49.35 $\pm$ 0.44	-0.20 $\pm$ 0.33	+ 0.33 $\pm$ 0.19	+ 0.48 $\pm$ 0.21
TriviaQA	FiD	67.74 $\pm$ 0.29	+ 0.40 $\pm$ 0.24	+ 0.89 $\pm$ 0.03	+ 0.92 $\pm$ 0.10
COLA	RoBERTa	60.41 $\pm$ 0.22	+ 0.09 $\pm$ 0.08	+ 1.57 $\pm$ 1.20	+ 1.41 $\pm$ 1.14
SST	RoBERTa	94.95 $\pm$ 0.13	-0.30 $\pm$ 0.27	-0.23 $\pm$ 0.40	+ 0.19 $\pm$ 0.14
MRPC	RoBERTa	89.14 $\pm$ 0.57	+ 0.08 $\pm$ 0.49	+0.73 $\pm$ 0.43	+ 0.81 $\pm$ 0.38
STSBB	RoBERTa	90.40 $\pm$ 0.02	+0.00 $\pm$ 0.05	+0.38 $\pm$ 0.17	+ 0.35 $\pm$ 0.16
QQP	RoBERTa	91.36 $\pm$ 0.07	+0.01 $\pm$ 0.06	+0.08 $\pm$ 0.07	+0.06 $\pm$ 0.08
MNLI	RoBERTa	87.41 $\pm$ 0.09	+0.08 $\pm$ 0.11	+0.39 $\pm$ 0.02	+0.35 $\pm$ 0.03
QNLI	RoBERTa	92.96 $\pm$ 0.06	-0.08 $\pm$ 0.11	+0.09 $\pm$ 0.01	+0.11 $\pm$ 0.06
RTE	RoBERTa	80.09 $\pm$ 0.23	-0.23 $\pm$ 0.20	+0.70 $\pm$ 0.65	-0.46 $\pm$ 0.12

(a)

Task	Model	Baseline	SWA	SAM	SWA+SAM
NPP: Proteins	SAGE	77.79 $\pm$ 0.18	-0.17 $\pm$ 0.22	-0.02 $\pm$ 0.13	-0.11 $\pm$ 0.15
	DGCN	85.42 $\pm$ 0.17	+ 0.11 $\pm$ 0.08	-0.14 $\pm$ 0.05	-0.08 $\pm$ 0.07
NPP: Products	SAGE	78.92 $\pm$ 0.08	+ 0.39 $\pm$ 0.10	+ 0.13 $\pm$ 0.08	+ 0.57 $\pm$ 0.03
	DGCN	73.88 $\pm$ 0.13	+ 0.44 $\pm$ 0.14	+ 0.08 $\pm$ 0.09	+ 0.53 $\pm$ 0.05
GPP: Code2	GCN	16.04 $\pm$ 0.09	+ 0.73 $\pm$ 0.11	+ 0.36 $\pm$ 0.08	+ 0.93 $\pm$ 0.15
	GIN	15.73 $\pm$ 0.11	+ 0.83 $\pm$ 0.11	+ 0.57 $\pm$ 0.09	+ 1.10 $\pm$ 0.09
GPP: Molpcba	GIN	28.10 $\pm$ 0.11	+ 0.40 $\pm$ 0.18	-0.33 $\pm$ 0.14	+ 0.33 $\pm$ 0.16
	DGCN	25.65 $\pm$ 0.13	+ 1.90 $\pm$ 0.20	-0.13 $\pm$ 0.18	+ 1.34 $\pm$ 0.12
LPP: Biokg	CP	84.06 $\pm$ 0.00	+ 0.07 $\pm$ 0.01	0.00 $\pm$ 0.03	+ 0.08 $\pm$ 0.02
	Complex	84.94 $\pm$ 0.01	+ 0.14 $\pm$ 0.01	-0.02 $\pm$ 0.01	+ 0.12 $\pm$ 0.02
LPP: Citation2	GCN	79.52 $\pm$ 0.41	-0.05 $\pm$ 0.52	+ 1.32 $\pm$ 0.06	+ 1.50 $\pm$ 0.13
	SAGE	81.95 $\pm$ 0.02	+ 1.15 $\pm$ 0.02	-0.31 $\pm$ 0.07	+ 0.86 $\pm$ 0.04

(b)

# Основные выводы:

1. Оптимизаторы поиска широкого оптимума могут уступать по ошибке другим оптимизаторам
2. SWA+SAM более устойчив и в 39/42 случаях показал результат лучше чем хотя бы один из SWA или SAM по отдельности
3. Рассмотренные оптимизаторы предлагают асимметричную отдачу: в худшем случае они снижают производительность на -0,30%, в лучшем - повышают на 2,60%.

# ИСТОЧНИКИ

SAM: <https://arxiv.org/abs/2010.01412>

SWA: <https://arxiv.org/pdf/1803.05407.pdf>

SAM+SWA вместе: <https://arxiv.org/abs/2202.00661>