

Momentum Residual Neural Networks¹

Докладчик - Алексей Цеховой, 181

Рецензент - Сибагатова Софья, 181

Практик-исследователь - Барановская Дарья, 181

Хакер - Андреев Никита, 181

¹ [M. E. Sander, et al. - "Momentum Residual Neural Networks", 2021.](#)

Цель:

- Хотим сократить память при обучении ResNet

Суть:

- Обратимая альтернатива ResNet: при обратном проходе активации вычисляются на ходу, при прямом не сохраняются

Дискретные:

- Обратимая связь активаций
- Сложность либо с переносом модели, либо с обращением

Непрерывные:

- ResNet как динамическая система, решаем ОДУ
- Сложность с переносом модели
- Численная неустойчивость

ResNet:

- $x_{n+1} = x_n + f(x_n, \theta_n)$

Momentum ResNet:

- $v_{n+1} = \gamma v_n + (1 - \gamma)f(x_n, \theta_n)$

- $x_{n+1} = x_n + v_{n+1}$

обратный проход:

- $x_n = x_{n+1} - v_{n+1}$

- $v_n = 1/\gamma(v_{n+1} - (1 - \gamma)f(x_n, \theta_n))$

Алгоритм²:

- 1: **Input:** Information buffer i , value c , ratio n/d
- 2: $i = i \times d$ ▷ make room for new digit
- 3: $i = i + (c \bmod d)$ ▷ store digit lost by division
- 4: $c = c \div d$ ▷ divide by denominator
- 5: $c = c \times n$ ▷ multiply by numerator
- 6: $c = c + (i \bmod n)$ ▷ add digit from buffer
- 7: $i = i \div n$ ▷ shorten information buffer
- 8: **return** updated buffer i , updated value c

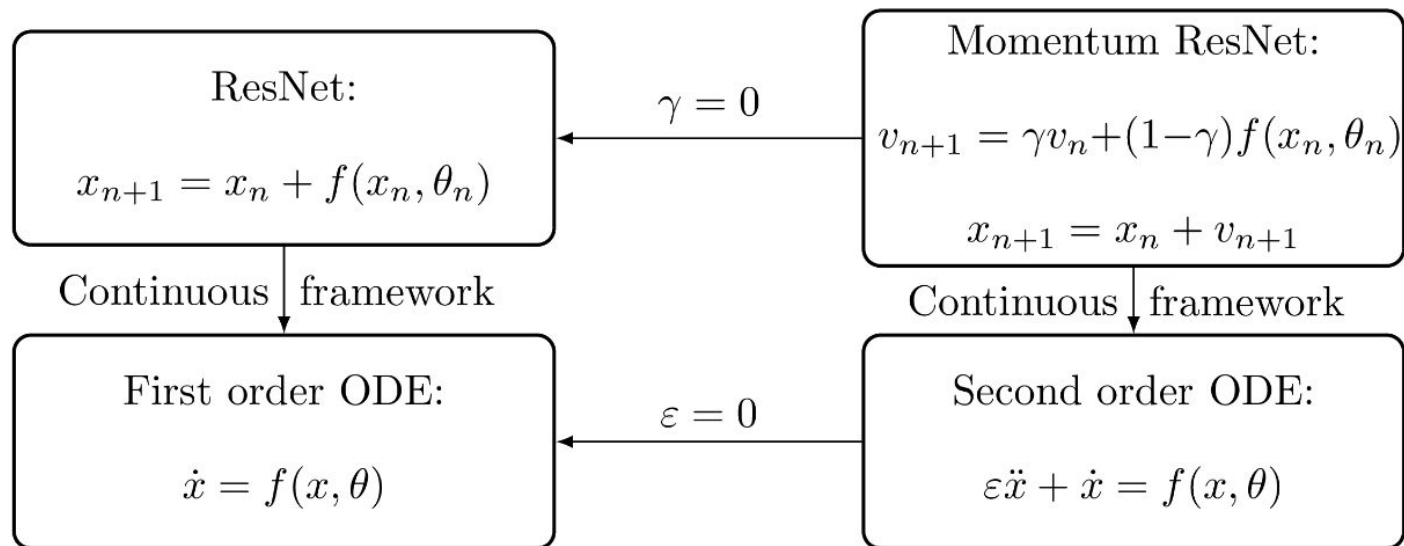
- Битов на операцию: $\log_2(1/\gamma)$

² [D. Maclaurin, et al. - "Gradient-based hyperparameter optimization through reversible learning", 2015.](#)

Обратимые остаточные архитектуры:

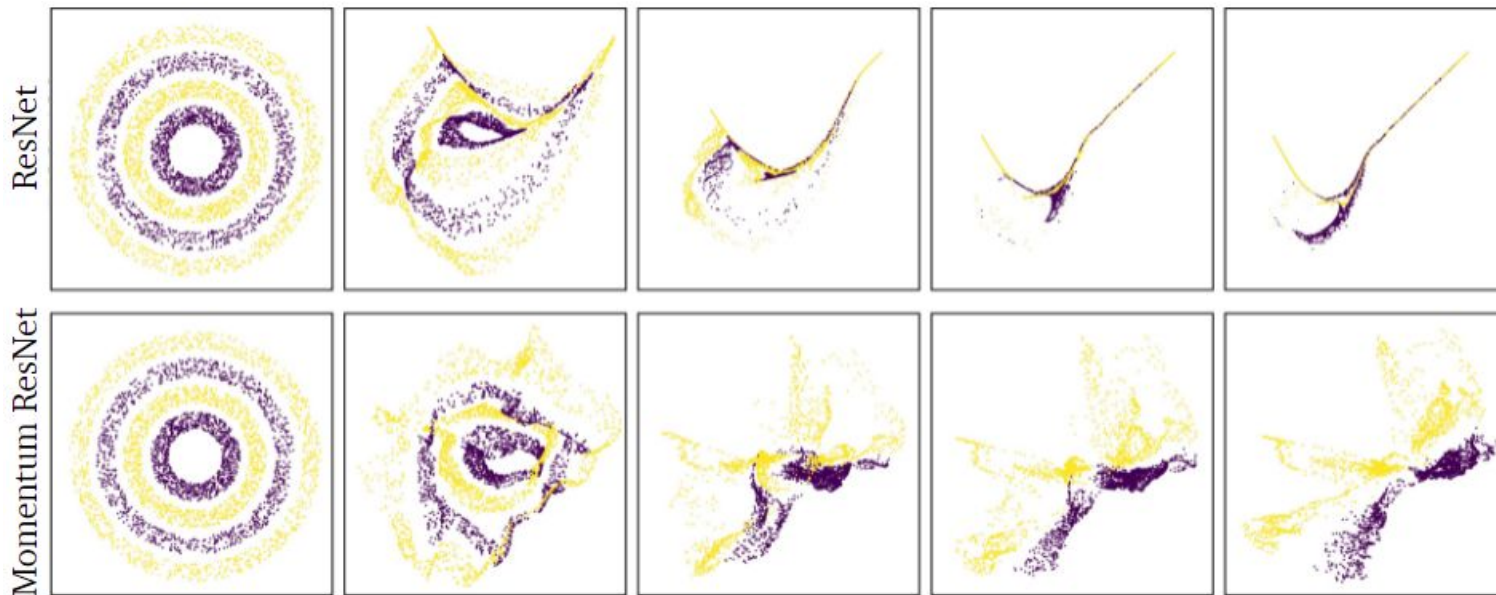
	<i>Neur.ODE</i>	<i>i-ResNet</i>	<i>i-RevNet</i>	<i>RevNet</i>	<i>Mom.Net</i>
Closed-form inversion	✓	✗	✓	✓	✓
Same parameters	✗	✓	✗	✗	✓
Unconstrained training	✓	✗	✓	✓	✓

- Обобщение ResNet (при моменте 0)
- Представляет ОДУ второго порядка, лучше ResNet - ОДУ первого порядка:



■ Опыт 1: лучшая возможность представления

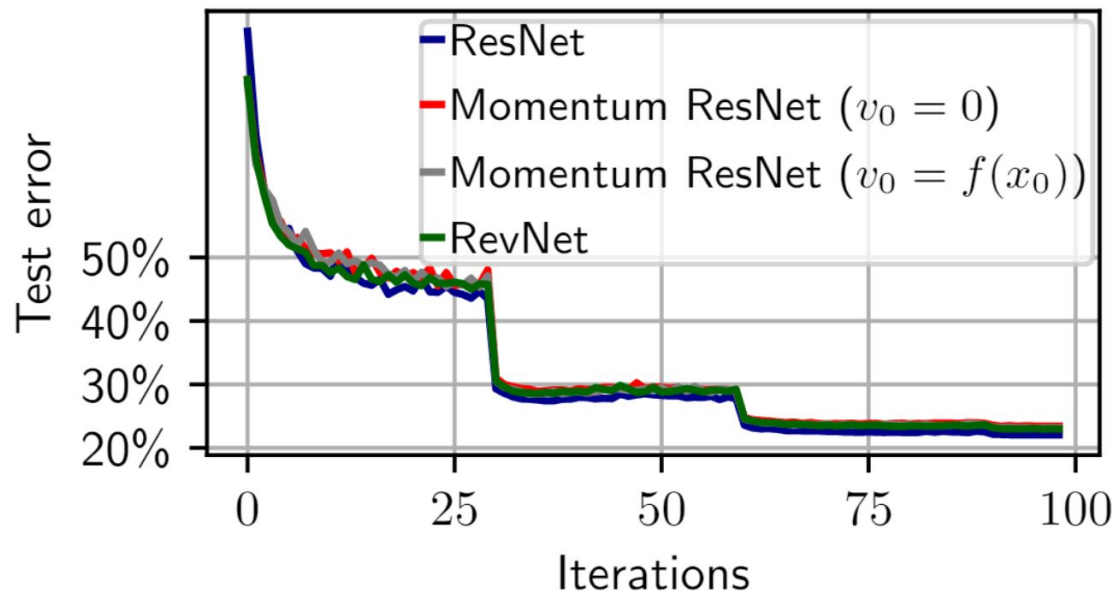
- $f(x, \theta) = W_2^T \tanh(W_1 x + b) : W_1, W_2 \in \mathbb{R}^{16 \times 2}, b \in \mathbb{R}^{16}$
- 15 слоёв
- Разделение колец (каждые 3 слоя):



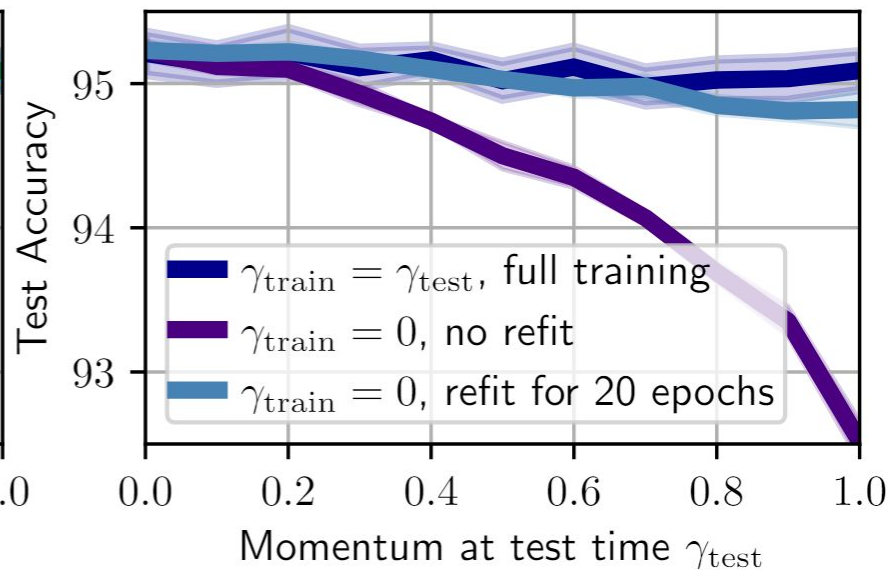
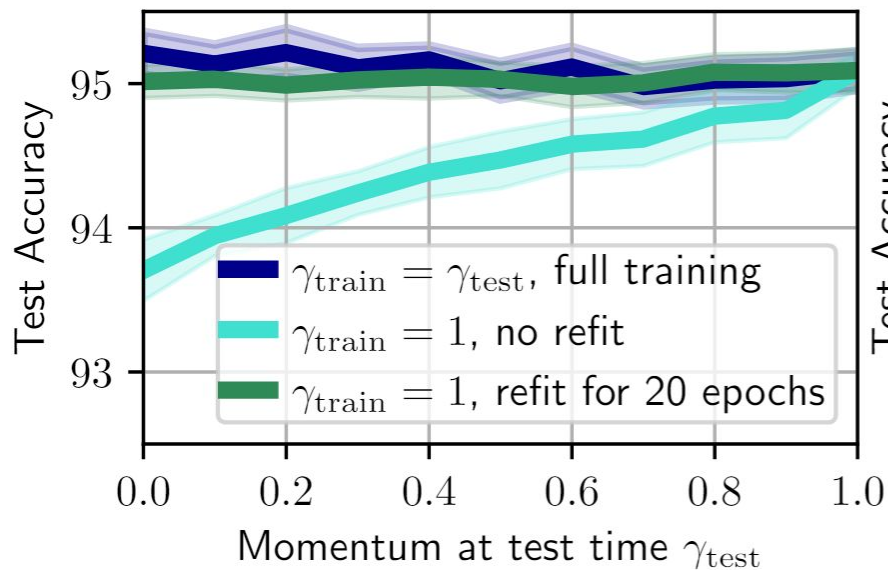
- ResNet-101 на CIFAR-10 и CIFAR-100
- $\gamma = 0.9$

Model	CIFAR-10	CIFAR-100
Momentum ResNet, $v_0 = 0$	95.1 ± 0.13	76.39 ± 0.18
Momentum ResNet, $v_0 = f(x_0)$	95.18 ± 0.06	76.38 ± 0.42
ResNet	95.15 ± 0.12	76.86 ± 0.25

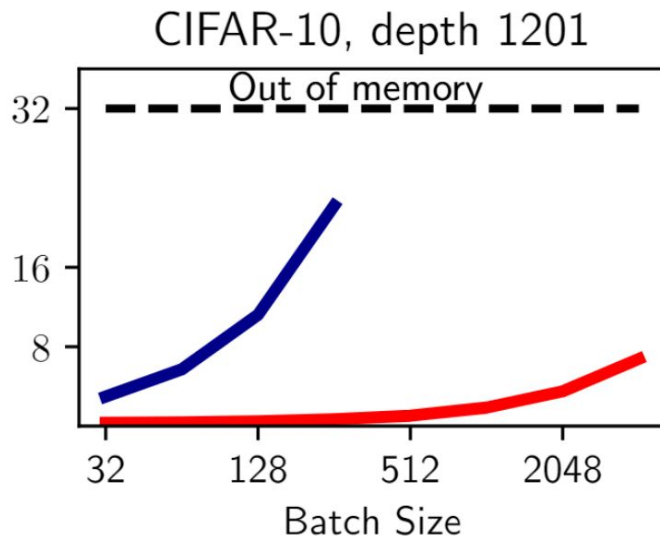
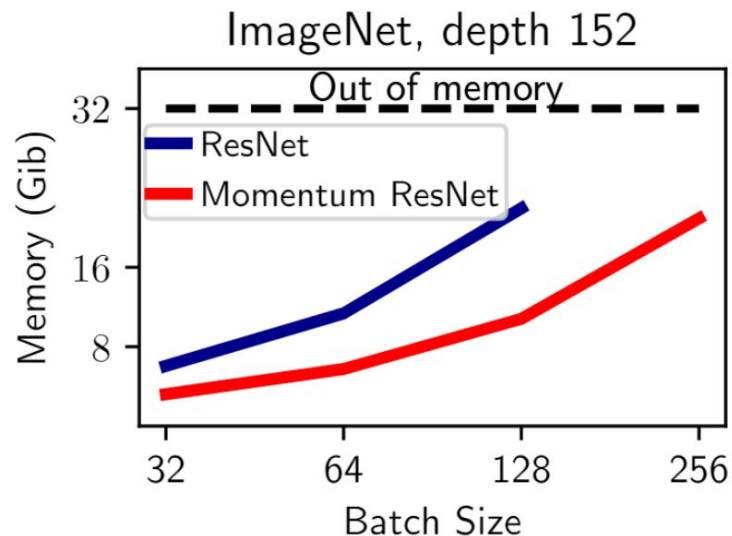
- ResNet-101 на ImageNet
- $\gamma = 0.9$
- ResNet: 22%, остальные: 23%



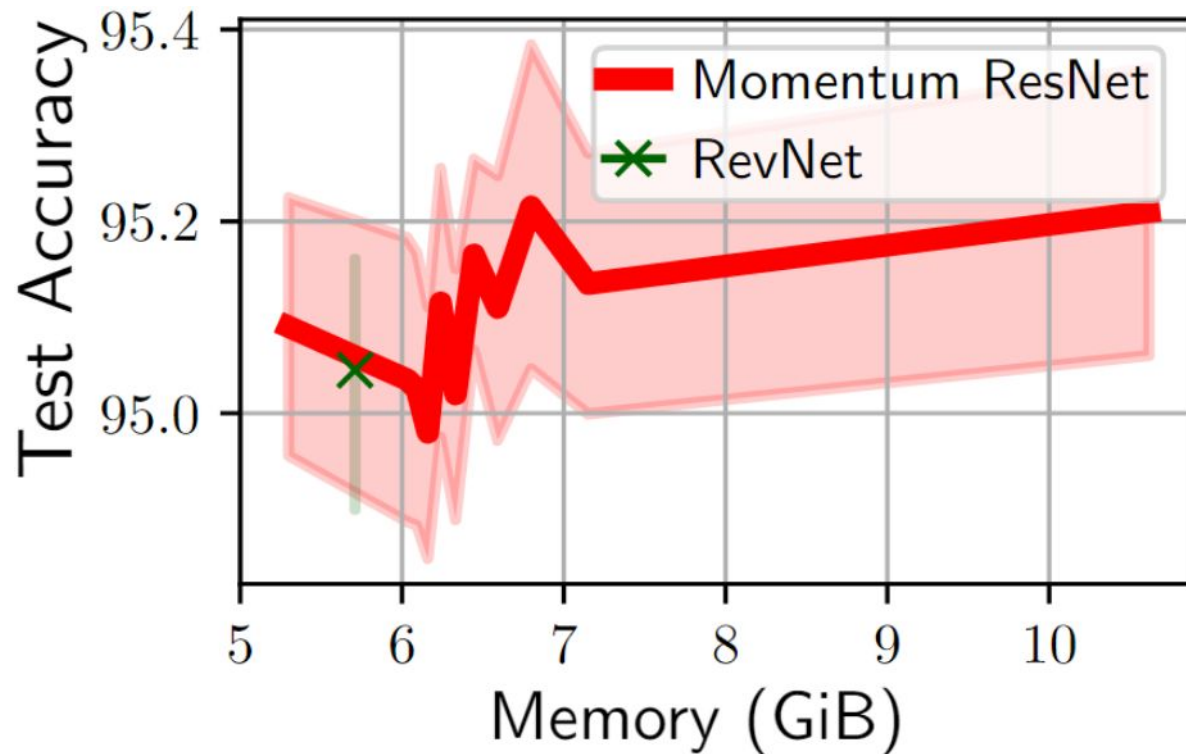
- ResNet-101 на CIFAR-10
- $\gamma = 0.9$



- ImageNet, глубина 152
- CIFAR-10, глубина 1201
- $\gamma = 0.9$

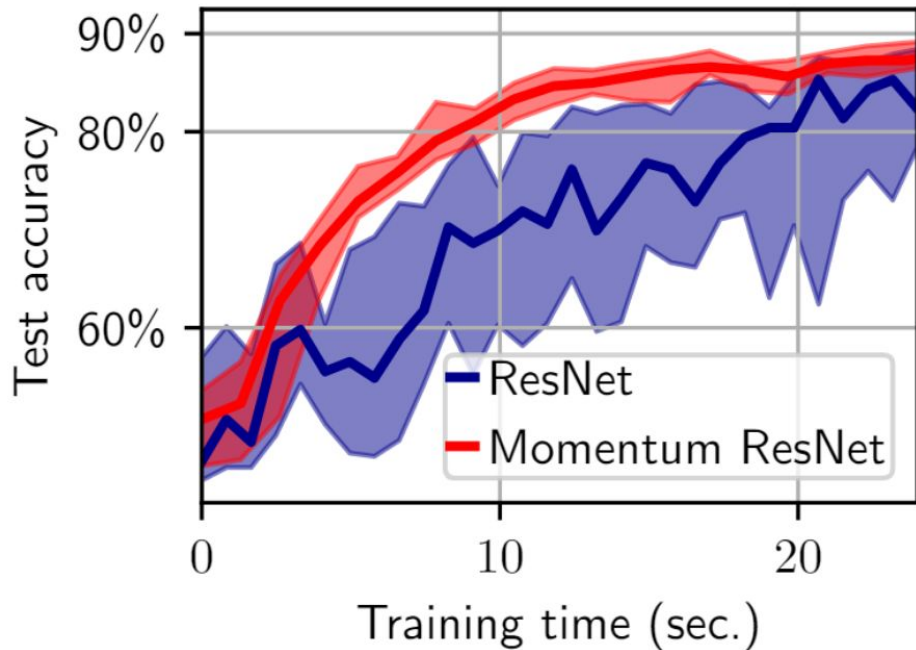


- ResNet-101 на CIFAR-10



- Обученный на Imagenet ResNet-152 (te=22%) перенесли в Momentum ResNet-152, $\gamma = 0.9$
- 1 эпоха дообучения: te=26.5%
- 5 эпох: te=23.5%

- Сети из опыта 7 на humenoptera (500x500px)
- Батчи по 2 и 4 изображения соответственно



- Значительно меньше памяти
- Сравнимое качество
- Простое построение из готовых ResNet
- Более широкий класс отображений

Momentum Residual Neural Networks

Michael E. Sander^{1 2} **Pierre Ablin**^{1 2} **Mathieu Blondel**³ **Gabriel Peyré**^{1 2}

ICML 2021 poster



Submission history

From: Michael E. Sander [[view email](#)]

[\[v1\]](#) Mon, 15 Feb 2021 22:24:52 UTC (2,224 KB)

[\[v2\]](#) Thu, 13 May 2021 12:29:54 UTC (4,555 KB)

[\[v3\]](#) Thu, 22 Jul 2021 08:18:05 UTC (4,554 KB)

Авторы



Michael E. Sander

Ecole Normale Supérieure

Machine Learning Applied Mathematics Optimal transport Differential equations



Pierre Ablin

CNRS

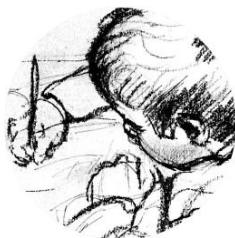
Optimization Machine Learning Statistics



Mathieu Blondel

Google

Machine Learning



Gabriel Peyré

CNRS, DMA, Ecole Normale Supérieure

Applied Mathematics Imaging sciences Machine learning

Deep Residual Learning for Image Recognition

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun

How Does Momentum Benefit Deep Neural Networks Architecture Design? A Few Case Studies

[PDF] [arxiv.org](#)

[B Wang](#), [H Xia](#), [T Nguyen](#), [S Osher](#) - arXiv preprint arXiv:2110.07034, 2021 - [arxiv.org](#)

We present and review an algorithmic and theoretical framework for improving neural network architecture design via momentum. As case studies, we consider how momentum can improve the architecture design for recurrent neural networks (RNNs), neural ordinary ...

☆ [🔗](#) [Цитировать](#) [Все версии статьи \(2\)](#) [🔗](#)

Interpolation and approximation via Momentum ResNets and Neural ODEs

[PDF] [arxiv.org](#)

[D Ruiz-Balet](#), [E Affili](#), [E Zuazua](#) - arXiv preprint arXiv:2110.08761, 2021 - [arxiv.org](#)

In this article, we explore the effects of memory terms in continuous-layer Deep Residual Networks by studying Neural ODEs (NODEs). We investigate two types of models. On one side, we consider the case of Residual Neural Networks with dependence on multiple ...

☆ [🔗](#) [Цитировать](#) [Все версии статьи \(6\)](#) [🔗](#)

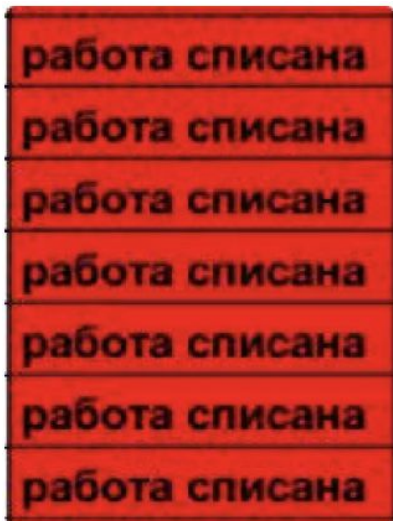
Efficient and Accurate Gradients for Neural SDEs

[PDF] [arxiv.org](#)

[P Kidger](#), [J Foster](#), [X Li](#), [T Lyons](#) - arXiv preprint arXiv:2105.13493, 2021 - [arxiv.org](#)

Neural SDEs combine many of the best qualities of both RNNs and SDEs, and as such are a natural choice for modelling many types of temporal dynamics. They offer memory efficiency, high-capacity function approximation, and strong priors on model space. Neural SDEs may ...

☆ [🔗](#) [Цитировать](#) [Цитируется: 1](#) [Похожие статьи](#) [Все версии статьи \(3\)](#) [🔗](#)



Плагнат

Alexander Chernov

вас предупредили



m-RevNet: Deep Reversible Neural Networks with Momentum

Duo Li[†] Shang-Hua Gao[‡]

[†]The Hong Kong University of Science and Technology [‡]Nankai University

duo.li@connect.ust.hk

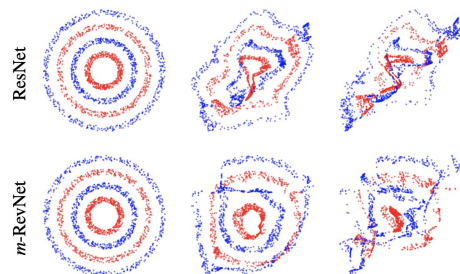
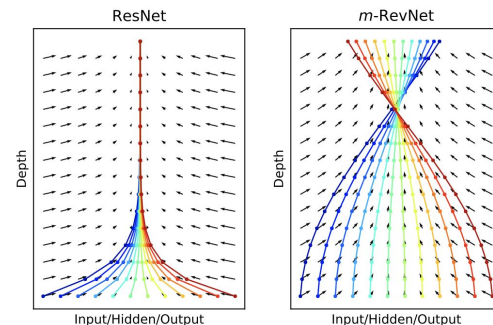
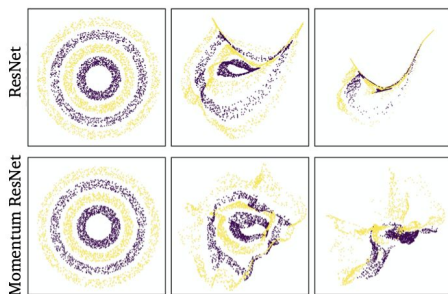
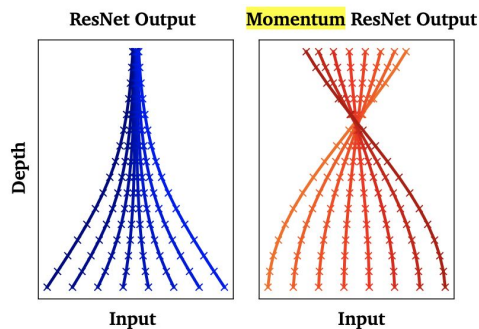
Submission history

From: Duo Li [[view email](#)]

[v1] Thu, 12 Aug 2021 17:14:32 UTC (4,231 KB)

[v2] Mon, 16 Aug 2021 13:04:04 UTC (0 KB)

- Только не скатывай
точь в точь
- Хорошо



- + Актуальность решаемой проблемы
 - + Применимые на практике результаты
 - + Большое количество обширных теоретических и экспериментальных исследований
-
- Недостаточно подробное описание некоторых экспериментов
 - В ключевых экспериментах о сравнении качества и потребляемой памяти Momentum ResNet и ResNet представлены результаты только для ResNet101 и его Momentum-аналога

- + Текст написан последовательно, подробно, доходчиво
- Большое количество орфографических и лингвистических ошибок, затрудняющих восприятие

3.2. Memory cost

Instead of storing the full data at each layer, we only need to store the bits lost at each multiplication by γ (cf. “invertibility”). For an architecture of depth k , this corresponds



invertibility

Top-1 classification error on ImageNet (single crop) for 4 different residual architectures of meters. Final test accuracy is 22% for the ResNet-101 and 23% for the 3 other invertible models. ie performance as a RevNet with the same number of parameters.



error

- + Предоставлен пакет для PyTorch, позволяющий получить Momentum-аналог любой модели ResNet (на самом деле также подходит для трансформеров)
- + Полное описание проводимых экспериментов
- Недостаточно подробное для воспроизведения описание эксперимента с облаками точек, но предоставлен код эксперимента, в котором есть все нужные гиперпараметры

- Предоставить результаты для ResNet различной глубины
- Провести сравнительные эксперименты времени обучения ResNet и Momentum ResNet
- Провести аналогичные для ResNet эксперименты по количеству памяти, качеству, возможности для fine-tuning моделей, содержащих residual блоки (например, трансформеры)

Оценка статьи по критериям НИПСа:

- Оценка: 9 из 10
- Уверенность: 5 из 5

Можно обращаться любые архитектуры с residual блоками. Для этого есть понятные функции и документация³.

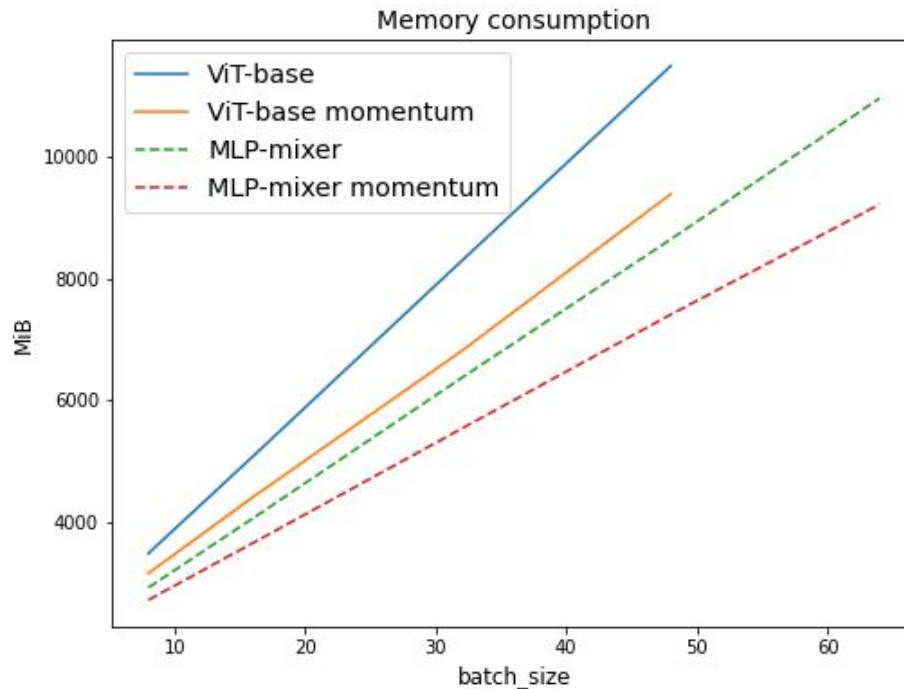
```
import torch
from momentumnet import transform_to_momentumnet
from torchvision.models import resnet101
resnet = resnet101(pretrained=True)
mresnet101 = transform_to_momentumnet(resnet, gamma=0.9, use_backprop=False)
```

```
import torch
from momentumnet import transform_to_momentumnet
transformer = torch.nn.Transformer(num_encoder_layers=6, num_decoder_layers=6)
mtransformer = transform_to_momentumnet(transformer, sub_layers=["encoder.layers", "decoder.layers"],
                                         gamma=0.9, use_backprop=False, keep_first_layer=False)
```

Таким образом можно экспериментировать с различными архитектурами и их моментум версиями.

³ <https://github.com/michaelsdr/momentumnet>

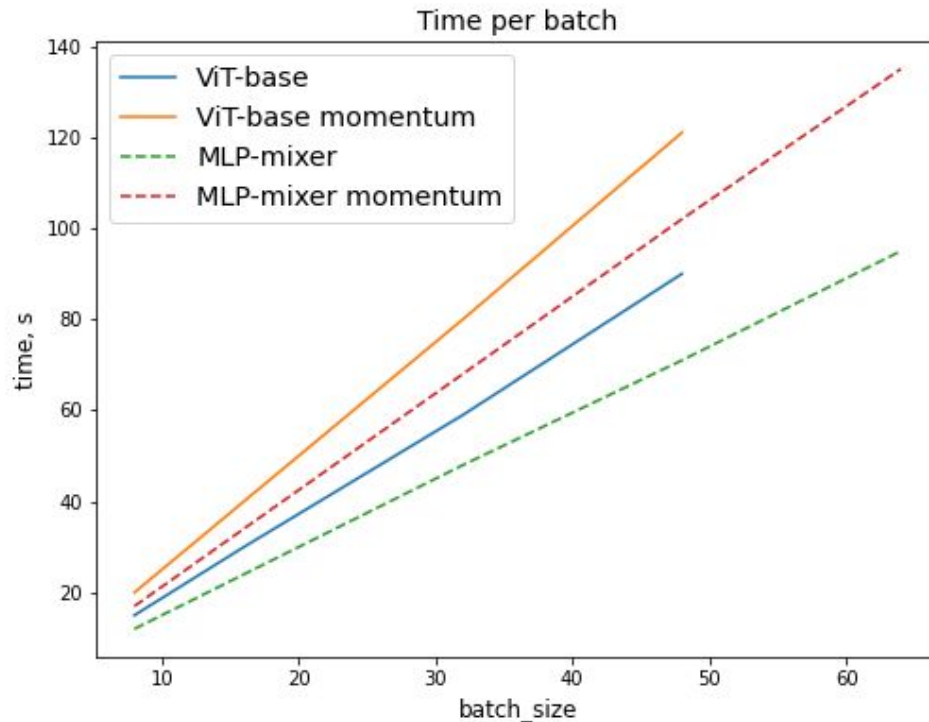
- Прогоняем модель на одном батче размера $bs \times 3 \times 224 \times 224$
- Эксперименты проводим с архитектурами Vision Transformer⁴ и MLP mixer⁵
- Batch_size 64 не влез для Vision Transformer
- Видим, что память растёт линейно, однако некоторая экономия все же наблюдается



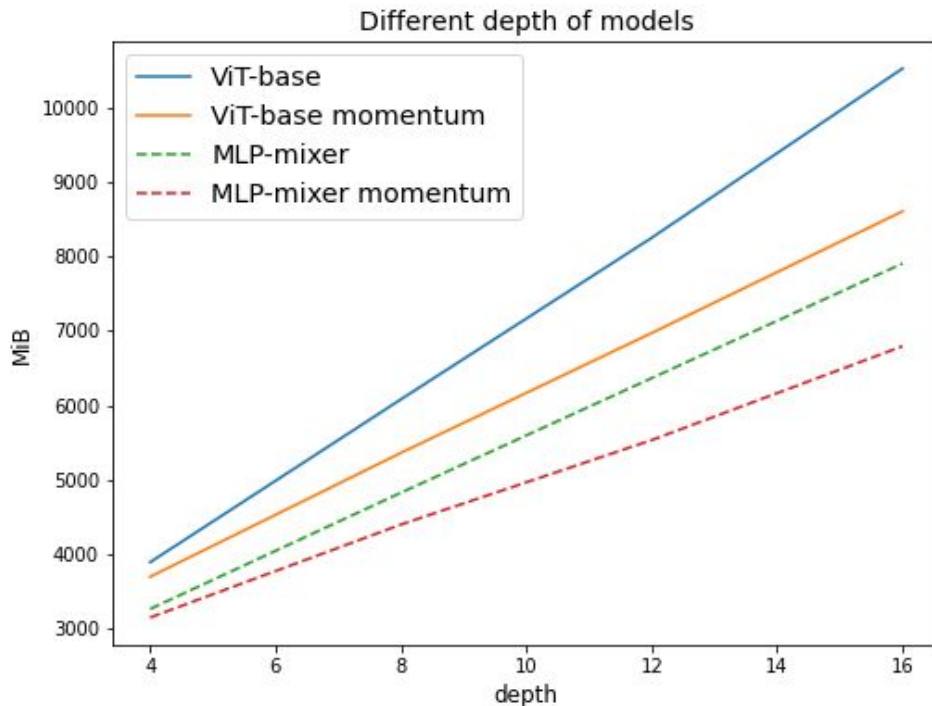
⁴ <https://arxiv.org/pdf/2010.11929.pdf>

⁵ <https://arxiv.org/pdf/2105.01601.pdf>

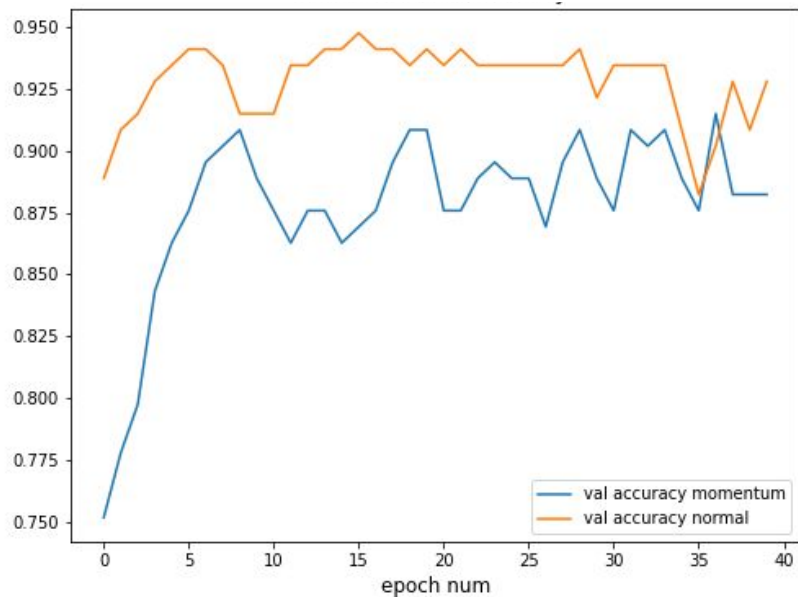
- Прогоняем модель на одном батче размера $bs \times 3 \times 224 \times 224$
- Эксперименты проводим с архитектурами Vision Transformer и MLP mixer
- Batch_size 64 не влез для Vision Transformer
- Видим, что время растет линейно, однако моментум версии проигрывают обычным моделям



- Прогоняем модель на одном батче размера 32 x 3 x 224 x 224
- Эксперименты проводим с архитектурами Vision Transformer и MLP mixer
- Видим, что при росте глубины модели (параметр depth) также есть эффект от применения моментум модификации



ResNet 18. Validation accuracy



ResNet 18. Validation loss

