

A ConvNet for the 2020s

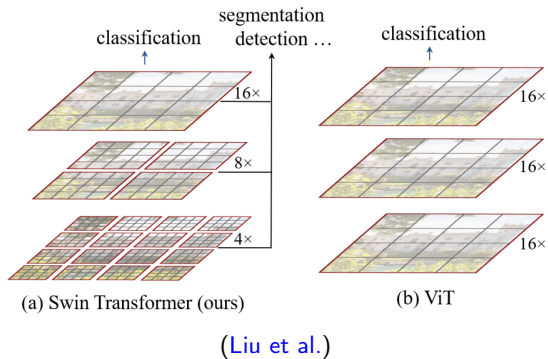
Шипилов Фома

ВШЭ ФКН ПМИ

5 октября 2022 г.

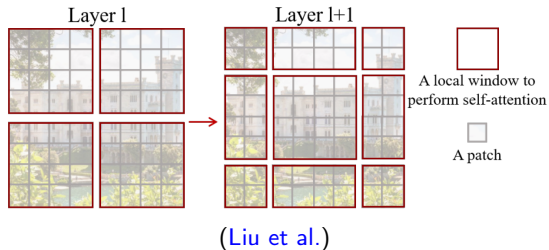
Swin Transformer: Ресурсы

- Иерархическая архитектура;



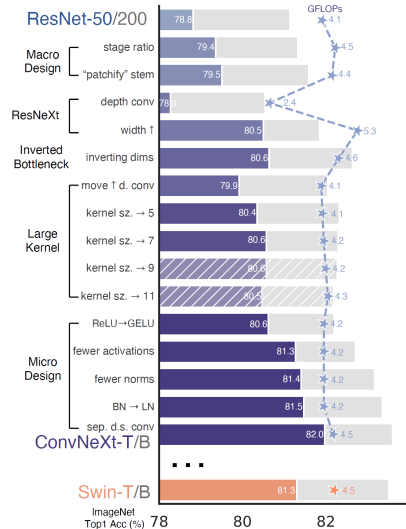
Swin Transformer: Рекап

- Иерархическая архитектура;
- Замена глобального self-attention на локальный windowed attention.



Roadmap

- 1 Бейзлайн: ResNet-50;
- 2 Пайплайн обучения;
- 3 Конфигурация блоков;
- 4 Патчизация;
- 5 ResNeXt;
- 6 Аналог windowed attention;
- 7 Функции активации;
- 8 Нормализации;
- 9 Даунсемплинг.



Что?

- 90 эпох → 300 эпох;

Image

ResNet-50



Mixup [47]



Cutout [3]



CutMix



Примеры аугментаций (Yun et al.)

Что?

- 90 эпох → 300 эпох;
- Adam → AdamW;

Image

ResNet-50



Mixup [47]



Cutout [3]



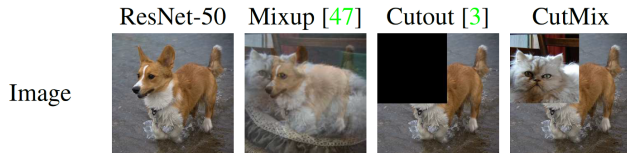
CutMix



Примеры аугментаций (Yun et al.)

Что?

- 90 эпох → 300 эпох;
- Adam → AdamW;
- Аугментации: Mixup, Cutmix, ...;



Примеры аугментаций (Yun et al.)

Что?

- 90 эпох → 300 эпох;
- Adam → AdamW;
- Аугментации: Mixup, Cutmix, ...;
- Регуляризации: Stochastic Depth, Label Smoothing.



Примеры аугментаций (Yun et al.)

Что?

- 90 эпох → 300 эпох;
- Adam → AdamW;
- Аугментации: Mixup, Cutmix, ...;
- Регуляризации: Stochastic Depth, Label Smoothing.

Зачем?

- Пайлайн как у Swin Transformer и DeiT



Примеры аугментаций (Yun et al.)

Что?

- $\text{Blocks}(3, 4, 6, 3) \rightarrow \text{Blocks}(3, 3, 9, 3)$

Зачем?

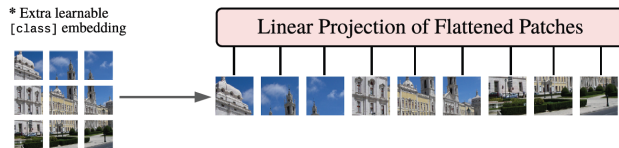
- Swin-T: $\text{Blocks}(1, 1, 3, 1)$;
- Объем вычислений на каждом уровне как у Swin-T.

Что?

- ResNet.in_conv: Conv2d(3, 64, kernel_size=7, stride=2) + MaxPool2d(2) \rightarrow Conv2d(3, 64, kernel_size=4, stride=4)

Зачем?

- Swin: токены — патчи 4×4 ;
- Более агрессивный даунсемплинг.



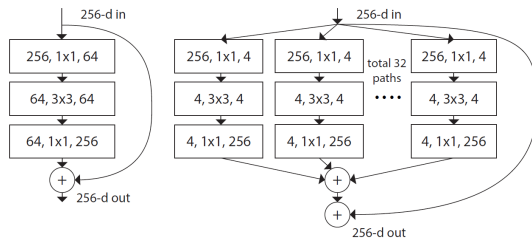
Патчизация ViT ([Dosovitskiy et al.](#))

Что?

- $\text{Conv2d}(64, 64, \text{kernel_size}=3) \rightarrow 96 \times \text{Conv2d}(96, 1, \text{kernel_size}=1) +$
 $+ 96 \times \text{Conv2d}(1, 1, \text{kernel_size}=3) +$
 $+ 96 \times \text{Conv2d}(1, 96, \text{kernel_size}=1) + \text{Sum}()$
 aka Depthwise(96, 96, hidden=96, kernel_size=3) свертка из MobileNet

Зачем?

- Разделение по-канальных и пространственных операций как в трансформерах;
- Количество каналов как в Swin-T.



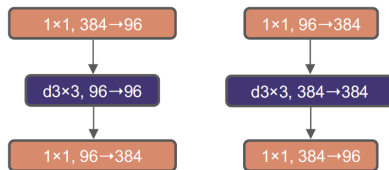
Слева: блок ResNet, справа: блок ResNeXt (Xie et al.)

Что?

- Depthwise(96, 96, hidden=96, kernel_size=3) → Depthwise(96, 96, hidden=384, kernel_size=3)
aka inverted bottleneck из MobileNetV2

Зачем?

- В 4 раза больше параметров в скрытом слое как в трансформерах



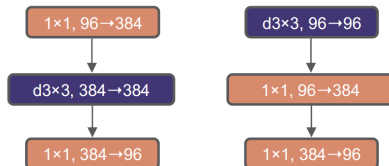
Слева: блок ResNeXt, справа: блок с inverted bottleneck.

Что?

- Depthwise(96, 96, hidden=384, kernel_size=3) →
→ $96 \times \text{Conv2d}(1, 1, \text{kernel_size}=3) + \text{Sum}() +$
+ $\text{Conv2d}(96, 384, \text{kernel_size}=1) +$
+ $\text{Conv2d}(384, 96, \text{kernel_size}=1)$
aka DepthwiseFirst(96, 96, hidden=384, kernel_size=3)

Зачем?

- В трансформерах attention применяется перед полносвязными слоями



Слева: блок с inverted bottleneck, **справа:** depthwise свертка сдвинута вверх.

Что?

- DepthwiseFirst(96, 96, hidden=384, kernel_size=3) →
→ DepthwiseFirst(96, 96, hidden=384, kernel_size=7)

Зачем?

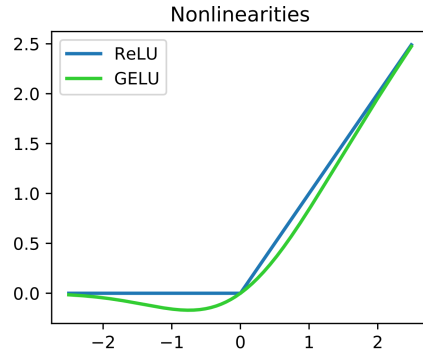
- Depthwise свертка аналогична windowed attention;
- Размер окна в Swin 7×7 .

Что?

- ReLU → GELU

Зачем?

- GELU используется в BERT, GPT-2, ViT, Swin



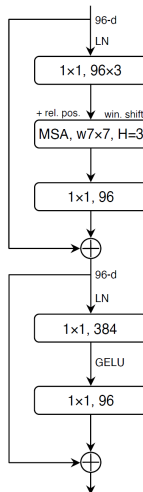
Что?

- Удаляем из блока все активации кроме одной

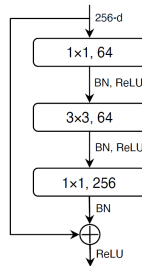
Зачем?

- Одна активация на блок как в трансформерах

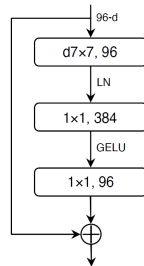
Swin Transformer Block



ResNet Block



ConvNeXt Block



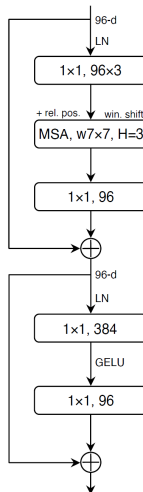
Что?

- Удаляем из блока все батч-нормализации кроме одной

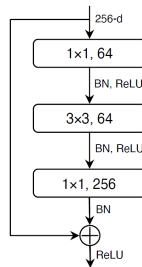
Зачем?

- Меньше нормализаций на блок как в трансформерах

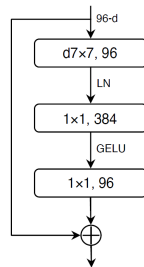
Swin Transformer Block



ResNet Block



ConvNeXt Block



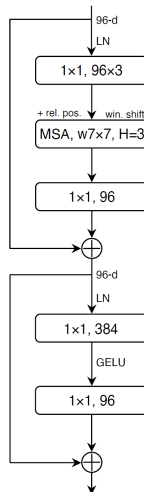
Что?

- Заменяем батч-нормализацию на LayerNorm

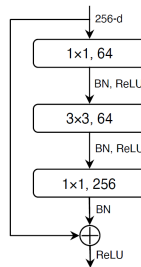
Зачем?

- LayerNorm используется в трансформерах

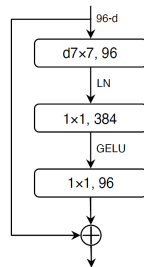
Swin Transformer Block



ResNet Block



ConvNeXt Block

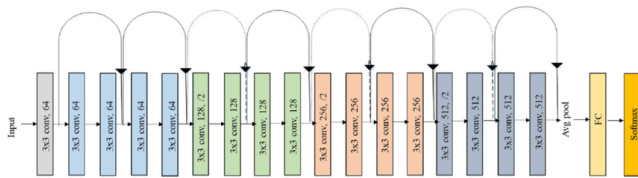


Что?

- ResNet.downsample: Conv2d(64, 128, kernel_size=3, stride=2) →
→ LayerNorm(96) +
+ Conv2d(96, 192, kernel_size=2, stride=2) +
+ LayerNorm(192) +
+ DepthwiseFirst(192, 192, hidden=768, kernel_size=7);
- Не делаем skip-connection через даунсемплинг.

Зачем?

- Стратегия даунсемплинга как в Swin



(He et al.)

Список источников

- ① A ConvNet for the 2020s
- ② Swin Transformer: Hierarchical Vision Transformer using Shifted Windows
- ③ CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features
- ④ An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale
- ⑤ Aggregated Residual Transformations for Deep Neural Networks
- ⑥ Deep Residual Learning for Image Recognition

ConvNeXt

от ResNet до SwinTransformer

Настя Городилова

БПМИ191

**Факультет Компьютерных Наук
НИУ ВШЭ**

4 октября 2022 г.

Когда и кем выпущена статья?

- ArXiv – январь 2022
- CVPR – июнь 2022

Авторы - Facebook AI Research

Откуда взялась идея: архитектура ResNeXt (CVPR 2017)

- 1 Aggregated residual transformations for deep neural networks.
CVPR, 2017 (ResNeXt)
- 2 Revisiting resnets: Improved training and scaling strategies.
NeurIPS, 2021
- 3 Resnet strikes back: An improved training procedure in timm.
CVPR, 2021

Прямые продолжения

- 1 VidConv: A modernized 2D ConvNet for Efficient Video Recognition, 2022
- 2 More ConvNets in the 2020s: Scaling up Kernels Beyond 51×51 using Sparsity, 2022

Плюсы статьи

- Конкурентоспособность ConvNeXt по отношению к Transformer-ам
- "Playthrough"

Минусы статьи

- Нет ветвления в поставленных экспериментах