

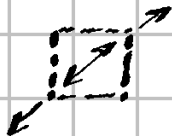
## EfficientNet V2 : Smaller Models and Faster Training.

Более быстрая работа, чем у предыдущих моделей, скорость обучения и лучшая эффективность параметров.

Используется комбинация поиска + масштабирование, чтобы оптимизировать скорость обучения и эффективность параметров.

Модель EfficientNet V2 меньше остальных, но обучается в 6.8 раз быстрее.

- Размер изображения  $\uparrow \Rightarrow$  скорость обучения  $\uparrow$   
 $\Rightarrow$  точность  $\downarrow$



Чтобы компенсировать  $\downarrow$  точности  $\Rightarrow$  авторы создали метод, который адаптивно увеличивает регуляризацию (например, увеличение данных) вместе с размером из-з.

## ① Введение:

Предлагаемая модель - EfficientNet. Исследование показало, что:

- 1) Обучение с очень большими размерами изображений происходит медленно
- 2) Свертки по глубине происходят медленно на ранних слоях
- 3) Одинаковое масштабирование на каждом этапе не является оптимальным.

Чтобы решить проблемы выше авторы создали пространство поиска с доп. операциями `Fused.MVConv` и применением `NAS` с учетом обучения и масштабирования для совместной оптимизации точности модели, ск. обучения и размера параметров. Такая модель обучается в 4 раза быстрее, при том, что она в 6.8 раз меньше по кол-ву параметров.

## ② EfficientNet V1

Повышение точности модели CNN связано с ↑ их размера и выч. сложности.

EfficientNet - семейство моделей, предложенных для баланса м/у точностью и эффективностью.

### ○ Проблема масштабирования НС.

#### • Традиционные методы масштабирования:

- 1) Масштабирование по глубине (Depthwise) - увеличение кол-ва слоев в сети.
- 2) Масштабирование по ширине (Widthwise) - увеличение кол-ва каналов в каждом слое.
- 3) Масштабирование по разрешению (Resolution Scaling) - увеличение разрешения входных изображений.

#### • Ограничения одностороннего масштабирования:

Масштабирование только по 1 из этих направлений может привести к неэффективному исп. ресурсов:

- 1) Только глубина: очень глубокие сети могут страдать от проблем с обучением, таких как исчезающий градиент.
- 2) Только ширина: широкие сети требуют больше памяти и выч. ресурсов.
- 3) Только разрешение: высокое разрешение изображений ↑ выч. сложность и может привести к переобучению.

⇒ надо найти оптимальный путь.

## o EfficientNet V1: Составное масштабирование (Compound Scaling).

### • Основная идея:

EfficientNet V1 вводит составное масштабирование, которое одновременно масштабирует глубину, ширину и разрешение изображения сети с помощью одного коэф. масштабирования  $\phi$ .

### • Формула составного масштабирования.

- Глубина:  $d = \alpha^q$  (depth)
- Ширина:  $w = \beta^q$  (width)
- Разрешение:  $r = \gamma^q$  (resolution)

Здесь  $\alpha, \beta, \gamma$  - коэф., определяющие насколько нужно масштабировать каждый параметр, а  $\phi$  - общий коэф. масштабирования, контролирующий общую нагрузку на выч. рес-ы.

### • Ограничения на коэф-ты:

Коэффициенты подбираются таким образом, чтобы удовлетворяли следующим условиям:

$$\begin{cases} \alpha \cdot \beta^2 \cdot \gamma \approx 2 \\ \alpha, \beta, \gamma \geq 1 \end{cases}$$

## • Базовая модель EfficientNet

### • Использование нейроавтоматического поиска архитектуры (NAS)

Авторы использовали метод AutoML MNAS для поиска оптимальной базовой архитектуры сети, названной EfficientNet-B0. Этот подход учитывает как точность модели, так и её выр. эф-ть.

### • Особенности базовой архитектуры

- MB Conv блоки: использованы модифицированные свёртки с расширением (Mobile Inverted Bottleneck Convolution)

- Swish активация: нелинейная ф-ция активации вместо ReLU, т.к. Swish менее подвержена "исч-ю граф."

Swish:  $f(x) = x \cdot \text{sigmoid}(x)$ , для более сложных сл-в.

ReLU:  $f(x) = \max(0, x)$ , для 1-ых слоёв.

- Сжатие и возбуждение (Squeeze-and-Excitation): механизмы внимания для усиления признаков.

## • Применение составного масштабирования

### • Создание семейства моделей EfficientNet:

Используя базовую EfficientNet-B0, применяя сост. масштабирование с разными значениями  $\mu$ , были созданы модели от EfficientNet-B1 до EfficientNet-B7.

- Пример масштабирования:  
EfficientNet-B1  $\rightarrow \mu=1$   
EfficientNet-B2  $\rightarrow \mu=2$   
...  
EfficientNet-B7  $\rightarrow \mu=7$ .

### ③ EfficientNetV2

- Главное отличие - оптимизация для скорости обучения.
- Прогрессивное обучение (Progressive Learning)  
Эта техника состоит в том, что модель обучается сначала на данных с низкой разрешением и постепенно увеличивает разрешение входов по мере того, как обучается на более поздних этапах. Это снижает вых. нагрузку на хар. Этапах и позволяет адаптироваться к сложным данным.
- Гибридные блоки (Fused MBConv)

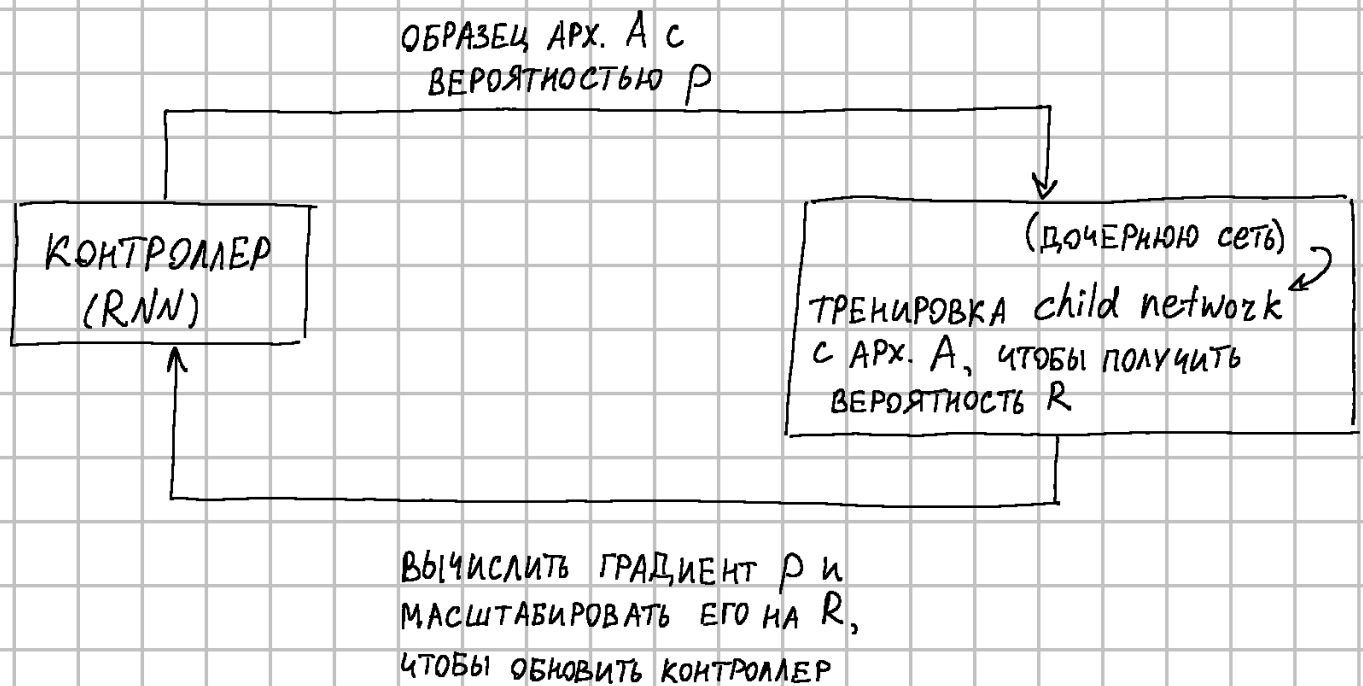
EfficientNetV2 использует новый тип блоков, который называется Fused MBConv. Эти блоки сочетают в себе традиционные сверточные операции и операции расширения (Expansion), что делает их более эффективными в плане вычислений по сравнению с классическими блоками MBConv, используемыми в 1-ой версии.

- Fused MBConv объединяет сверточные и расширенные операции в один шаг, что ↓ вых. затраты и ↑ произв-во.
- Fused MBConv эффективна на ранних этапах обучения (слои 1-3).

Quoc V. Le:

## 2016: Neural Architecture Search with Reinforcement Learning (NAS RL)

Основная идея заключается в том, что есть контроллер. Выбирается архитектура и производится оценка во время обучения, а также происходит проверка на иск. (валидационном) наборе. Полученная точность выдает на вознаграждение контроллера, которое представляет собой точность валидации  $R$ .



В итоге можно будет получить архитектуры с высокой точностью валидации.

## 2017: Learning Transferable Architectures for Scalable Image Recognition (LTAS IR)

Здесь говорится о переносе обучаемых архитектур для масштабируемого распознавания изображений.

Авторы создали область поиска, заметив, что большинство созданных вручную архитектур (например, мод. сети, RESNET) сводится к созданию ячейки и её последующему многократному повторению.

⇒ авторы решили разработать "ячейку", чтобы повторять её  $N$  раз.

$$\boxed{\bullet} \times N$$

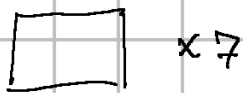
## 2018: MnasNet: Platform-aware Neural Architecture Search for Mobile.

Авторы увеличили обл. поиска:

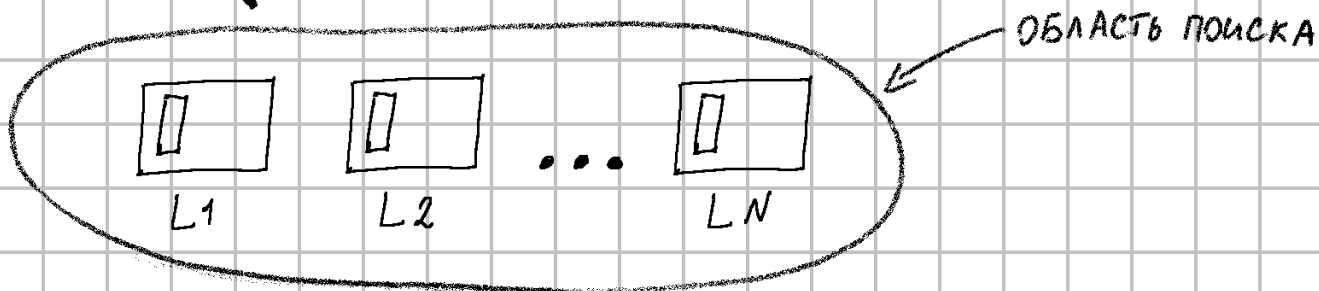




Теперь есть несколько блоков (пусть их будет 7)



Выбираем 1-ый слой в каждом из блоков, а затем повторяем его несколько раз, это также явл. гиперпараметры, кот. мы ищем.

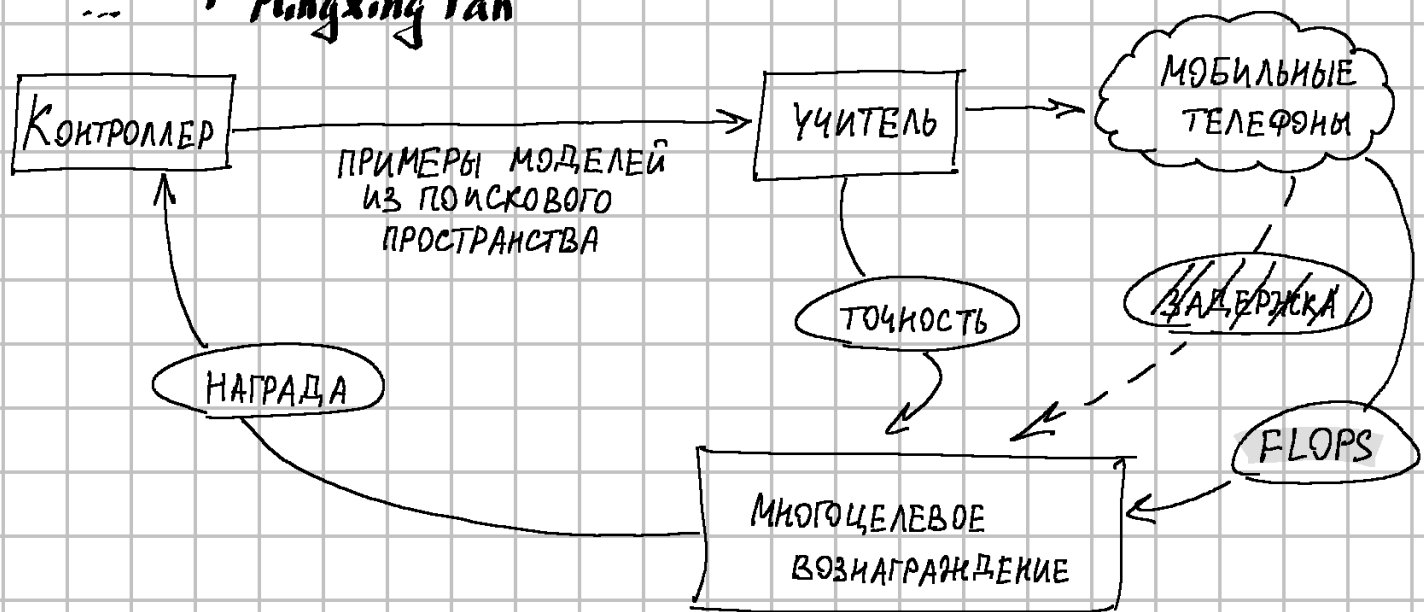


Также авторы включили задержку в конечную ф-цию вознаграждения  $\Rightarrow$  учитывается не только точность на проверочном наборе данных, но и задержка.

\* Это сделал статью подходящей для конца архитектур, кот. можно использовать на мод. центр-вах и др. периферийных центр-вах. Процесс поиска тот же, изм. только ф-ция вознаграждения. Тут не исп-ся ан. усиления, вместо него исп. PPO от OpenAI.

# 2019: EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks.

... + Mingxing Tan



Авторы взяли базовую архитектуру V0 и добавили в неё параметры.

2020: Pandemonium ...

2021: EfficientNet V2.