



Group of  
Horribly  
Optimistic  
Statisticians



CV SEMINAR

WE HAVE TO GO DEEPER!

09.01.2024 Computer Vision Seminar 23/24



**GHOST**

Group of Horribly Optimistic Statisticians



# Agenda

1. Inception - przypomnienie
2. Zanikający gradient
3. Połączenia “skrótowe”
4. ResNet
5. Tutorial transfer learning w PyTorch





**GHOST**

Group of Horribly Optimistic Statisticians



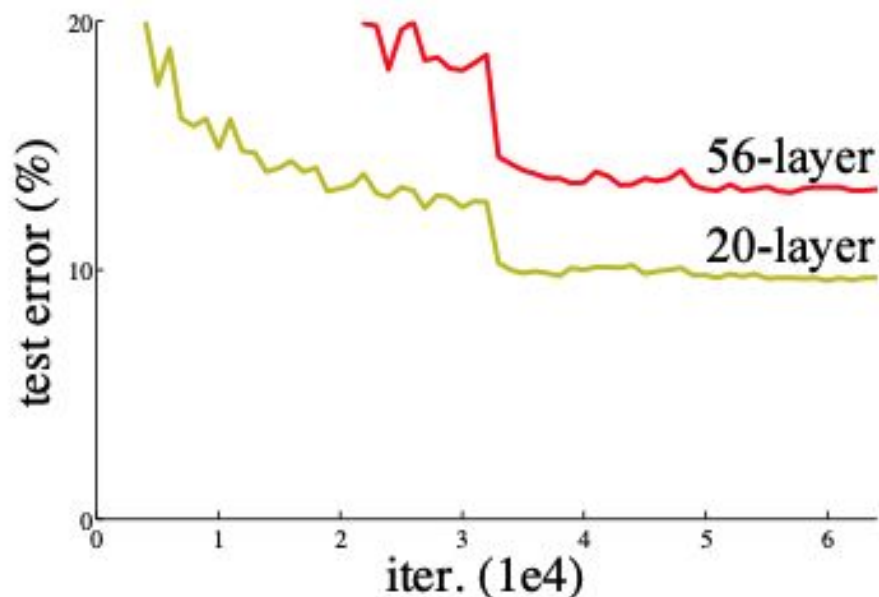
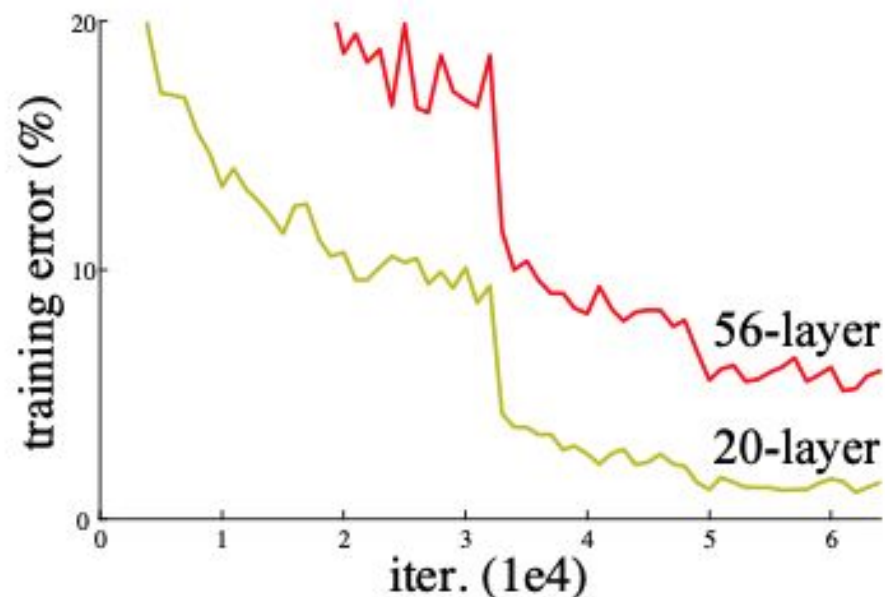


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.



**GHOST**

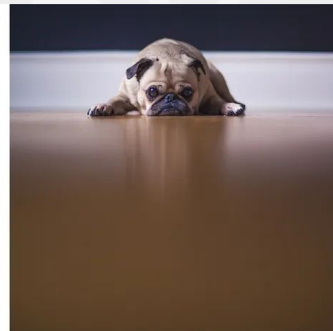
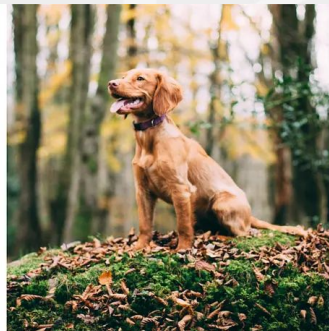
Group of Horribly Optimistic Statisticians



# Inception - motywacja

- Jak dobrać rozmiar filtra? (dla cech różnego rozmiaru)
- Bardzo głębokim sieciom grozi przeuczenie i zanik gradientu
- Układanie wielu warstw konwolucyjnych jest kosztowne obliczeniowo

<https://towardsdatascience.com/a-simple-guide-to-the-versions-of-the-inception-network-7fc52b863202>



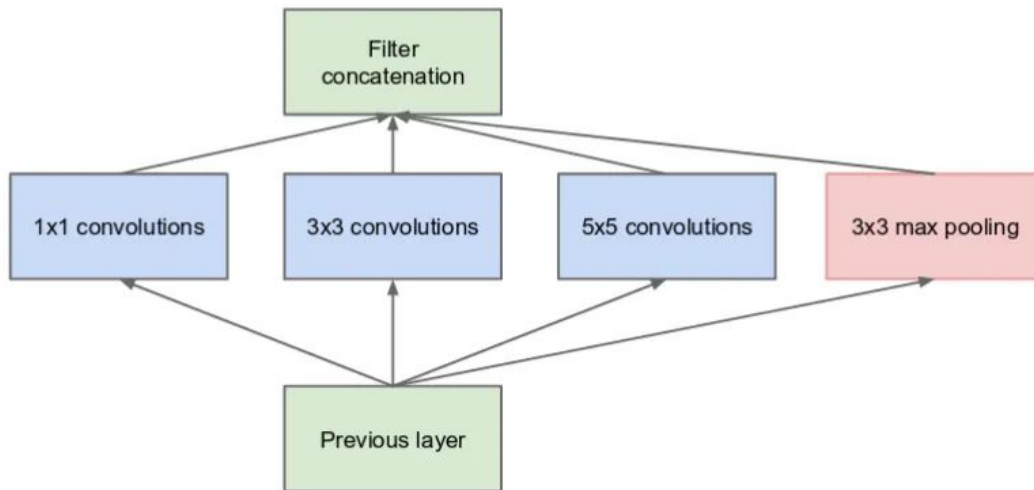


**GHOST**

Group of Horribly Optimistic Statisticians



## Rozwiązanie - różne wielkości filtrów na tym samym poziomie



(a) Inception module, naïve version

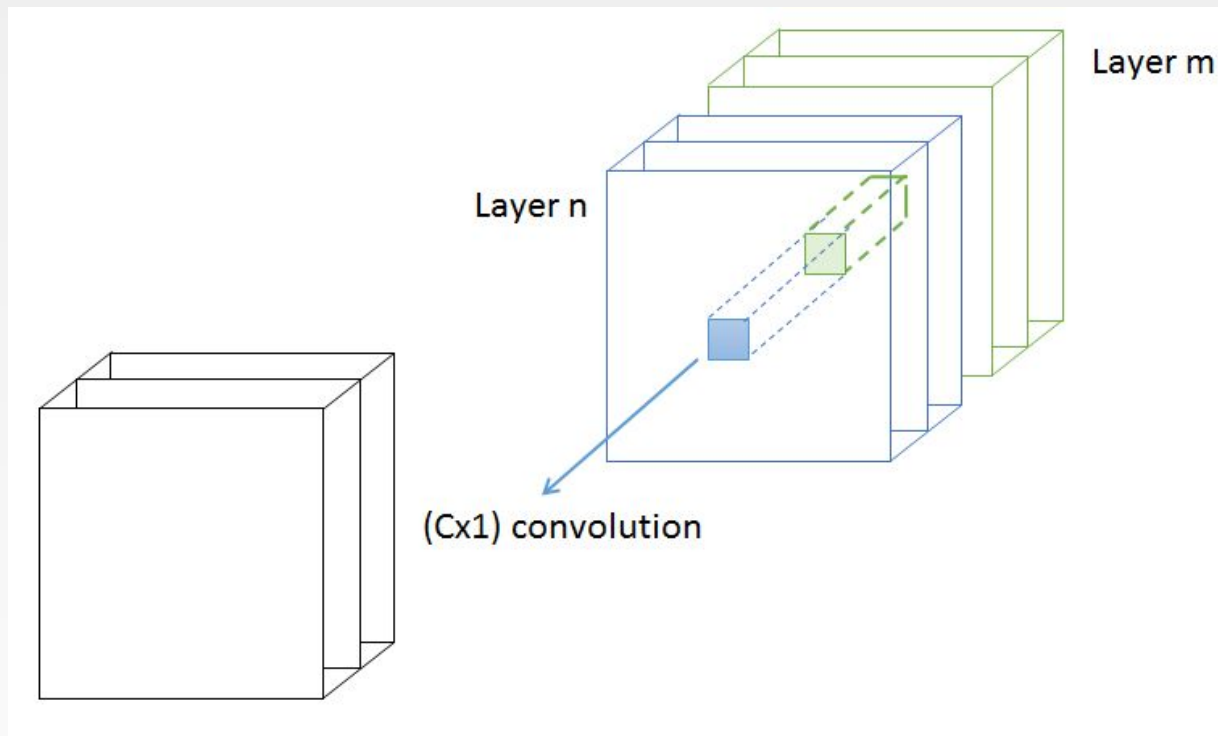


**GHOST**

Group of Horribly Optimistic Statisticians



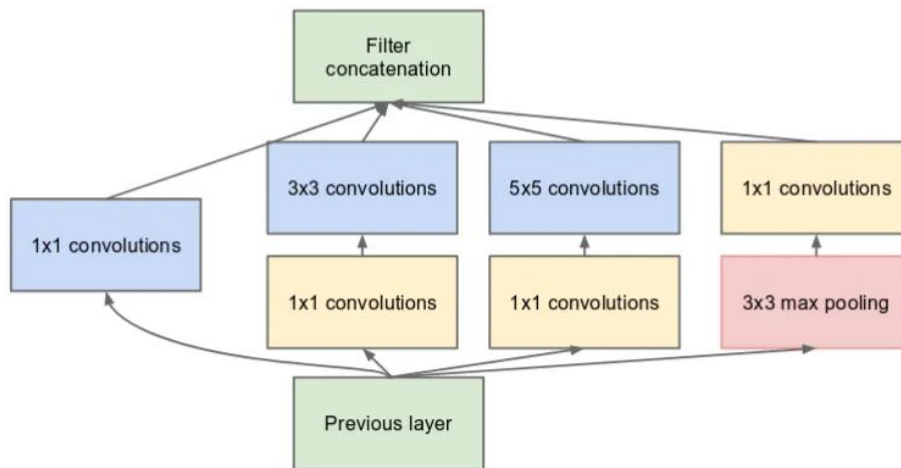
## Konwolucja 1x1





# Redukcja wymiarów

- Aby obniżyć koszt obliczeniowy, autorzy dodali konwolucję 1x1 (redukcja wymiarów).
- Konwolucja 1x1 jest znacznie mniej złożona obliczeniowo niż konwolucja 3x3 czy 5x5



(b) Inception module with dimension reductions



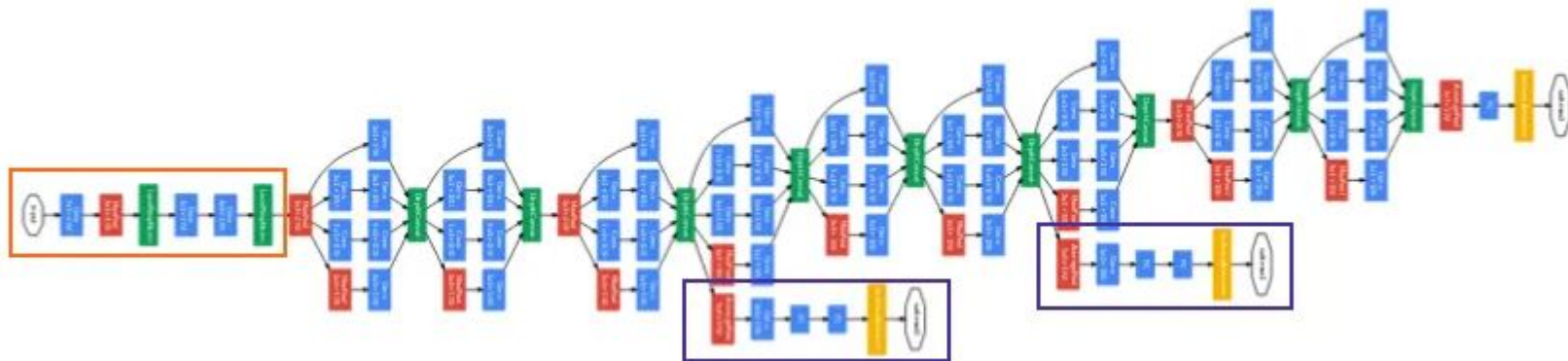


**GHOST**

Group of Horribly Optimistic Statisticians



# GoogLeNet





**GHOST**

Group of Horribly Optimistic Statisticians



# Zanik gradientu

- Gradient (wektor pochodnych) jest używany przy aktualizacji wag modelu
- Jeśli ułożymy kolejno kilka warstw o określonych funkcjach aktywacji, może dojść do zaniku
- Zbyt mały gradient uniemożliwia stabilne uczenie sieci

<https://towardsdatascience.com/the-vanishing-gradient-problem-69bf08b15484>

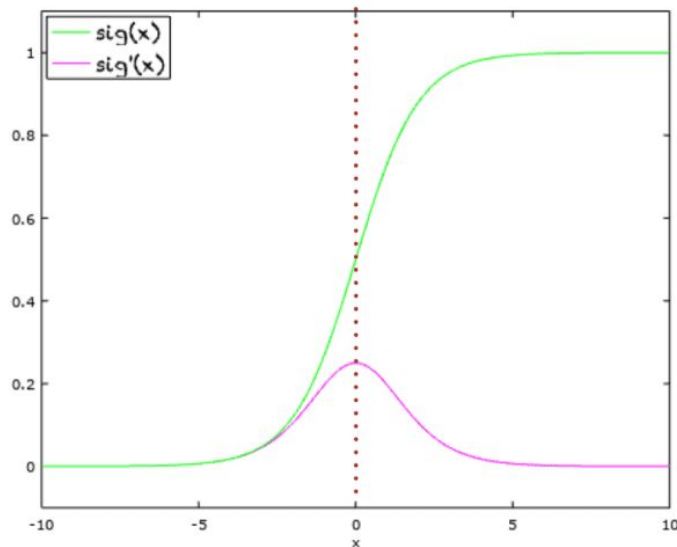


**GHOST**

Group of Horribly Optimistic Statisticians



# Sigmoida i jej pochodna - przykład



Plot of  $\sigma(x)$  and its derivate  $\sigma'(x)$

Domain:  $(-\infty, +\infty)$

Range:  $(0, +1)$

$\sigma(0) = 0.5$

Other properties

$$\sigma(x) = 1 - \sigma(-x)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$



**GHOST**

Group of Horribly Optimistic Statisticians



# Sigmoida i jej pochodna - przykład

- Duża zmiana wartości - mała zmiana pochodnej
- W propagacji wstecznej wartości pochodnej są mnożone przez siebie (chain rule, różniczkowanie funkcji złożonej)
- Jeśli kilka kolejnych warstw używa aktywacji sigmoidą, mnożymy kilka bardzo małych wartości, co powoduje dążenie gradientu do 0

<https://towardsdatascience.com/the-vanishing-gradient-problem-69bf08b15484>



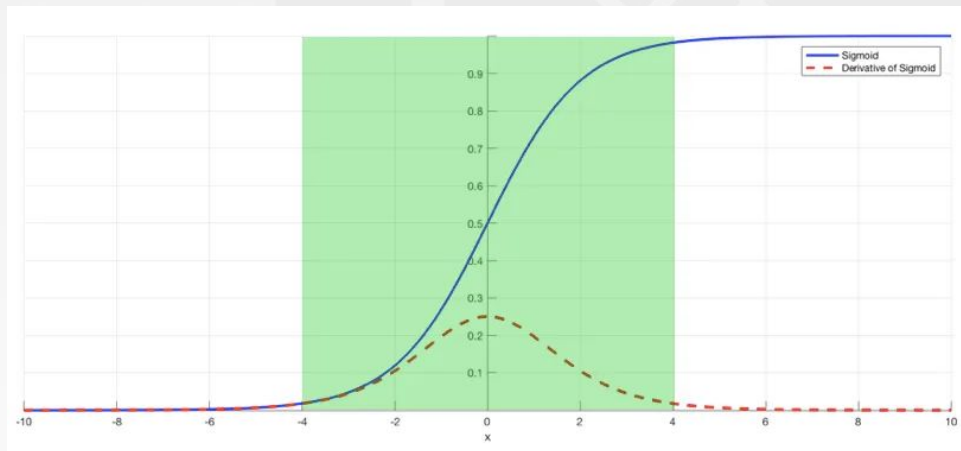
**GHOST**

Group of Horribly Optimistic Statisticians



# Zanikający gradient - rozwiązania

- Inne funkcje aktywacji, np. ReLU
- Batch normalization
- Połączenia rezydualne



<https://towardsdatascience.com/the-vanishing-gradient-problem-69bf08b15484>

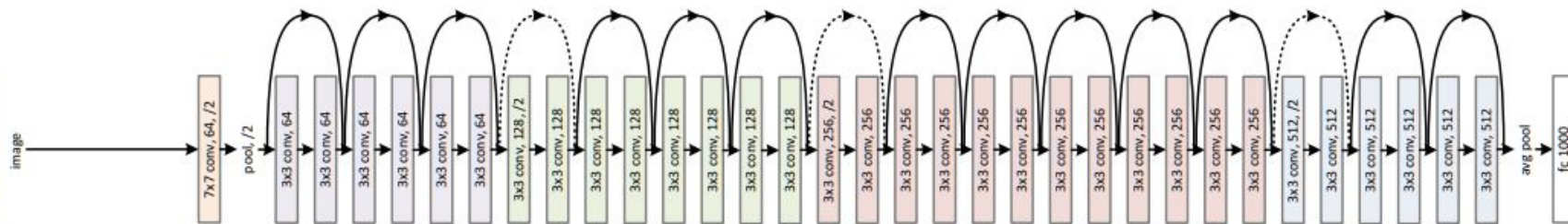


# GHOST

Group of Horribly Optimistic Statisticians



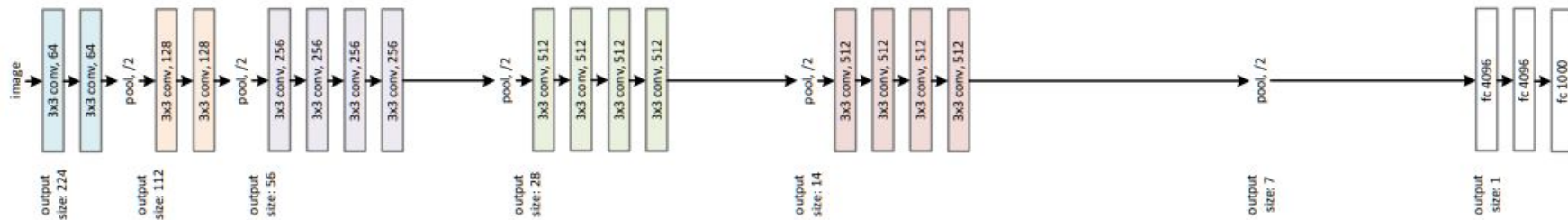
34-layer residual



34-layer plain



VGG-19



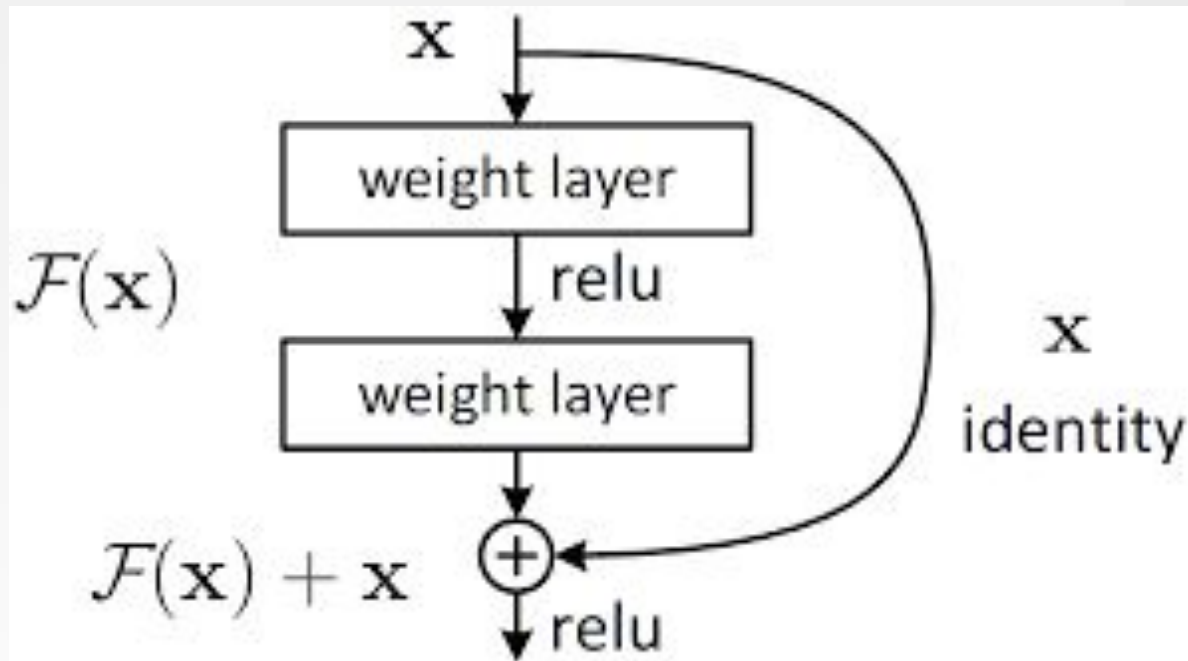


**GHOST**

Group of Horribly Optimistic Statisticians



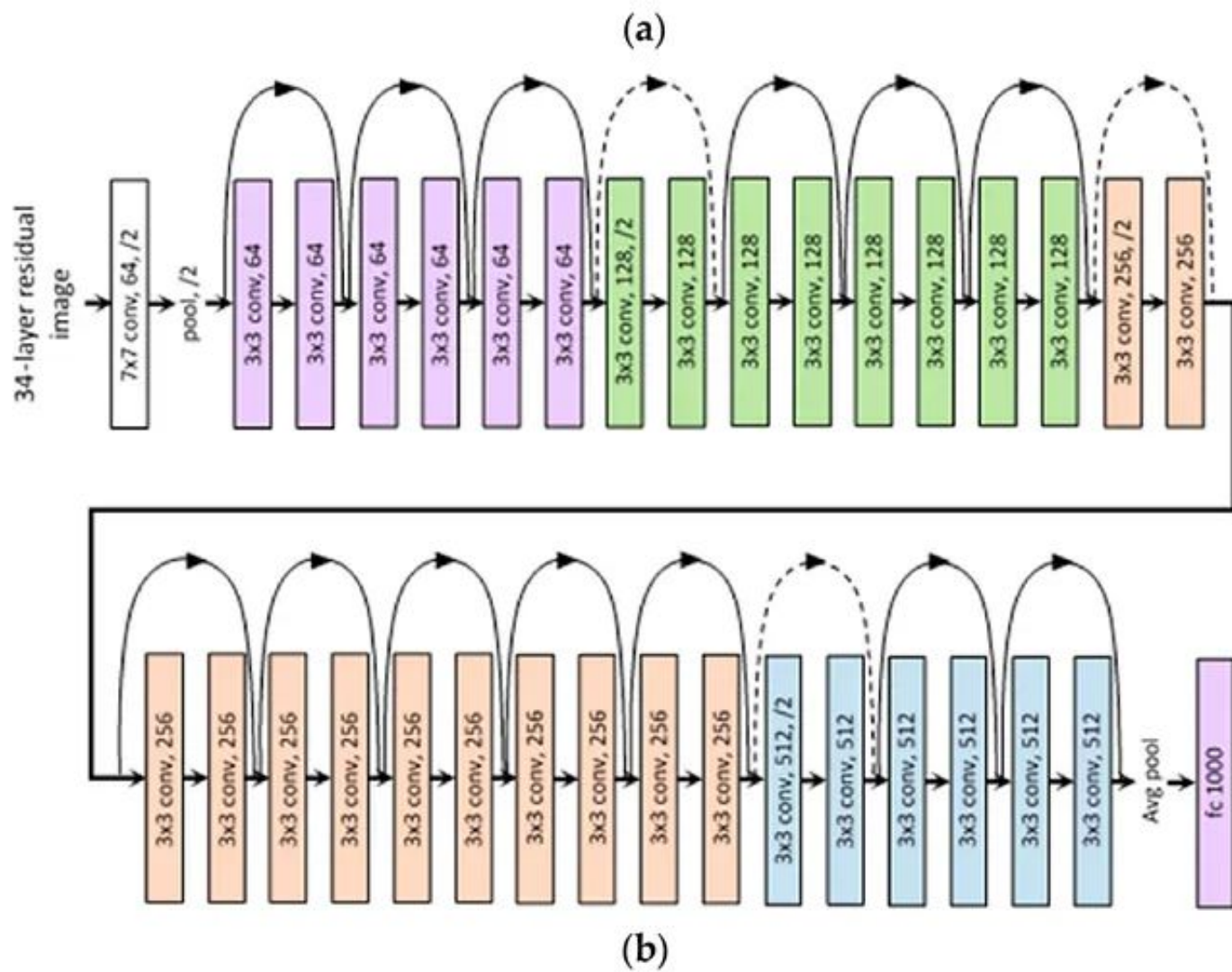
# Blok ResNet - połączenia “skrótowe”



In this paper, we address the degradation problem by introducing a *deep residual learning* framework. Instead of hoping each few stacked layers directly fit a desired underlying mapping, we explicitly let these layers fit a residual mapping. Formally, denoting the desired underlying mapping as  $\mathcal{H}(x)$ , we let the stacked nonlinear layers fit another mapping of  $\mathcal{F}(x) := \mathcal{H}(x) - x$ . The original mapping is recast into  $\mathcal{F}(x) + x$ . We hypothesize that it is easier to optimize the residual mapping than to optimize the original, unreferenced mapping. To the extreme, if an identity mapping were optimal, it would be easier to push the residual to zero than to fit an identity mapping by a stack of nonlinear layers.

The formulation of  $\mathcal{F}(x) + x$  can be realized by feedforward neural networks with “shortcut connections” (Fig. 2). Shortcut connections [2, 34, 49] are those skipping one or more layers. In our case, the shortcut connections simply perform *identity* mapping, and their outputs are added to the outputs of the stacked layers (Fig. 2). Identity shortcut connections add neither extra parameter nor computational complexity. The entire network can still be trained end-to-end by SGD with backpropagation, and can be easily implemented using common libraries (e.g., Caffe [19]) without modifying the solvers.







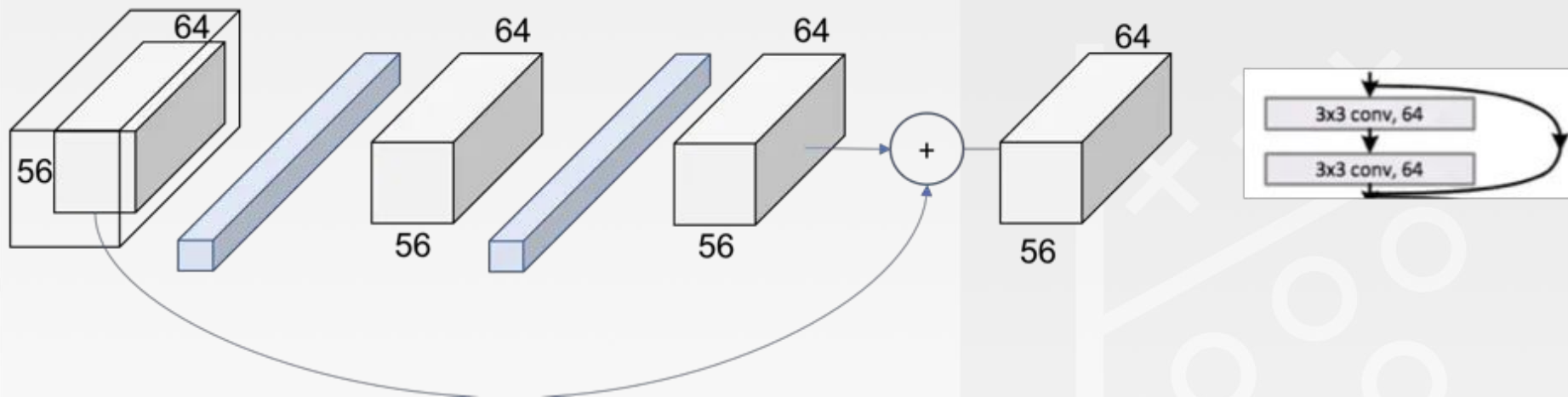


**GHOST**

Group of Horribly Optimistic Statisticians



# Blok ResNet



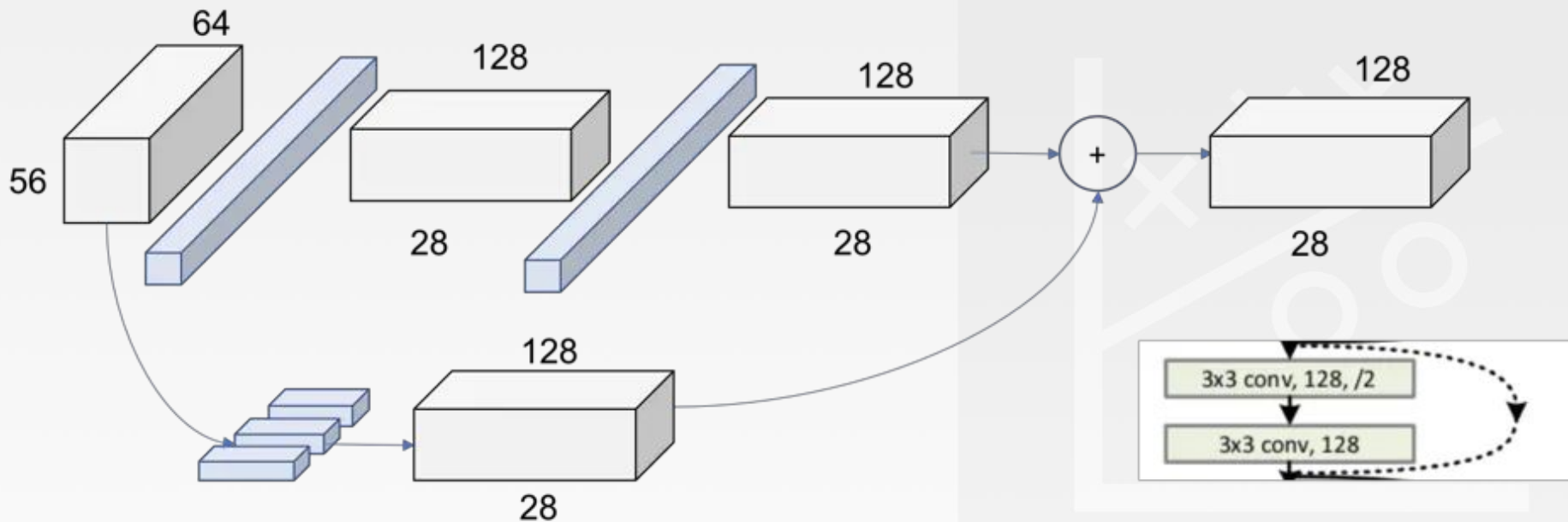


**GHOST**

Group of Horribly Optimistic Statisticians



# Blok ResNet - zwiększanie głębokości





## GHOST

Group of Horribly Optimistic Statisticians



model	top-1 err.	top-5 err.
VGG-16 [41]	28.07	9.33
GoogLeNet [44]	-	9.15
PReLU-net [13]	24.27	7.38
plain-34	28.54	10.02
ResNet-34 A	25.03	7.76
ResNet-34 B	24.52	7.46
ResNet-34 C	24.19	7.40
ResNet-50	22.85	6.71
ResNet-101	21.75	6.05
ResNet-152	<b>21.43</b>	<b>5.71</b>

Table 3. Error rates (% , **10-crop** testing) on ImageNet validation. VGG-16 is based on our test. ResNet-50/101/152 are of option B that only uses projections for increasing dimensions.

The Top-1 error is the proportion of the time the classifier does not provide the highest score to the correct class. The Top-5 error rate is the percentage of times the classifier failed to include the proper class among its top five guesses.



# GHOST

Group of Horribly Optimistic Statisticians



Number of Layers	Number of Parameters
ResNet 18	11.174M
ResNet 34	21.282M
ResNet 50	23.521M
ResNet 101	42.513M
ResNet 152	58.157M

Table 1. ResNets architectures for ImageNet



**GHOST**

Group of Horribly Optimistic Statisticians



# ResNet - podsumowanie

- Sieci rezydualne rozwiązują problem zanikającego gradientu, ponieważ gradient może swobodnie przepływać przez połączenia skrótowe
- Dzięki zastosowaniu połączeń skrótowych, sieci mogą być znacznie głębsze, co z kolei przekłada się na lepsze wyniki



**GHOST**

Group of Horribly Optimistic Statisticians

# Materialy

- TRANSFER LEARNING FOR COMPUTER VISION TUTORIAL

<https://colab.research.google.com/drive/1c9LNEwzSuSQrvFKzDBJ4yLiNgAVyDsV?usp=sharing>

[https://pytorch.org/tutorials/beginner/transfer\\_learning\\_tutorial.html](https://pytorch.org/tutorials/beginner/transfer_learning_tutorial.html)

- Understanding and visualizing ResNets

<https://towardsdatascience.com/understanding-and-visualizing-resnets-442284831be8>

