# Deep  Residual  Networks

su
2018/1116
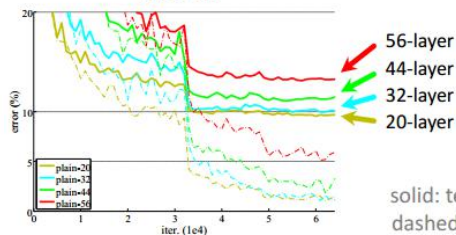
# Simply stacking layers?



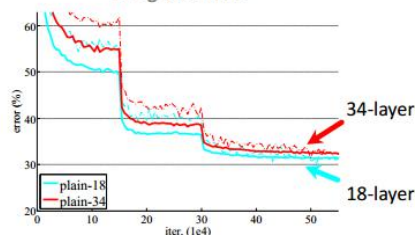CIFAR-10

train error (%) — 56-layer, 20-layer

test error (%) — 56-layer, 20-layer

- *Plain* nets: stacking 3x3 conv layers…
- 56-layer net has **higher training error** and test error than 20-layer net

CIFAR-10 — 56-layer, 44-layer, 32-layer, 20-layer

solid: test/val
dashed: train

ImageNet-1000 — 34-layer, 18-layer

- "Overly deep" plain nets have **higher training error**
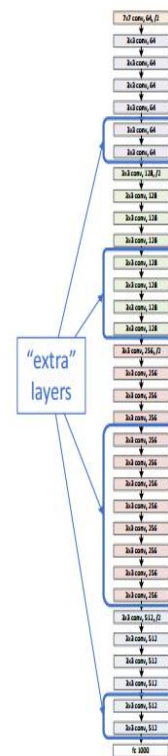- A general phenomenon, observed in many datasets

a shallower model (18 layers)

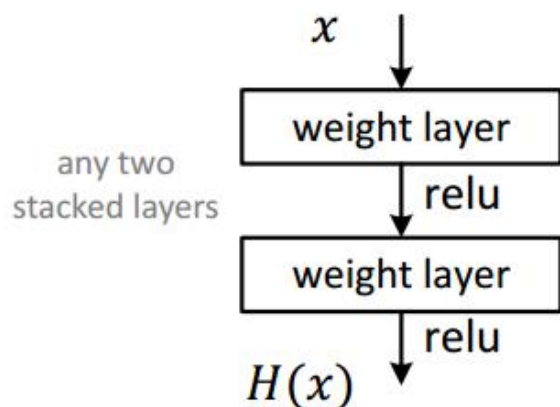a deeper counterpart (34 layers)

"extra" layers

- A deeper model should not have **higher training error**

- A solution *by construction*:
  - original layers: copied from a learned shallower model
  - extra layers: set as identity
  - at least the same training error

- Optimization difficulties: solvers cannot find the solution when going deeper…

# Deep Residual Learning

- Plaint net

$x$

weight layer

↓ relu

weight layer

↓ relu

$H(x)$

any two stacked layers

$H(x)$ is any desired mapping,

hope the 2 weight layers fit $H(x)$

- Residual net

$x$

weight layer

↓ relu

$F(x)$

weight layer

identity

$x$

$H(x) = F(x) + x$ ⊕

↓ relu

$H(x)$ is any desired mapping,

~~hope the 2 weight layers fit $H(x)$~~

hope the 2 weight layers fit $F(x)$
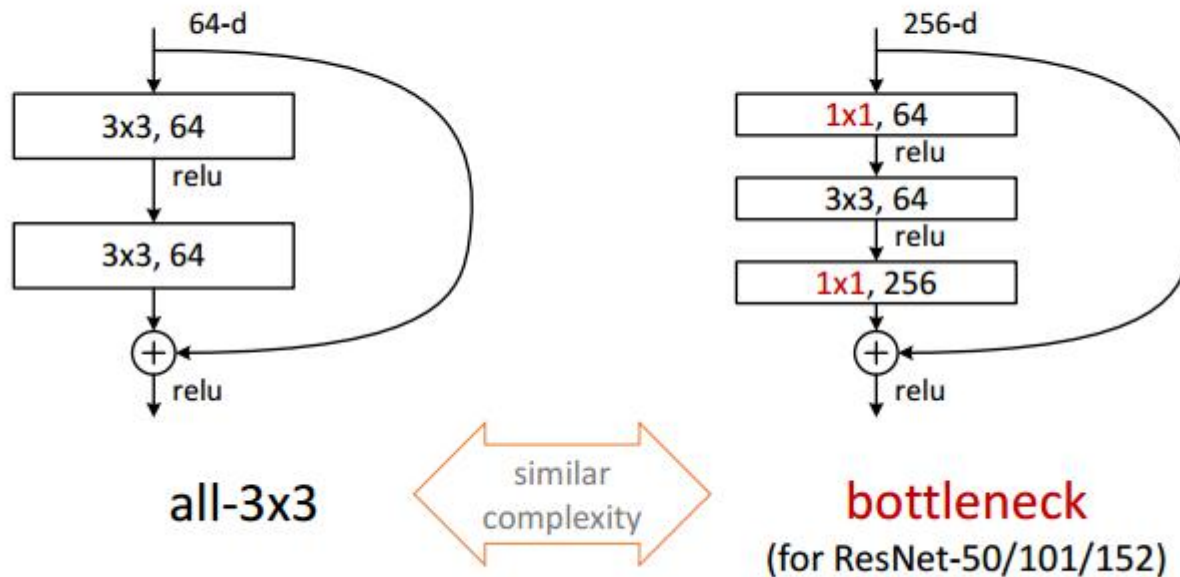
let $H(x) = F(x) + x$
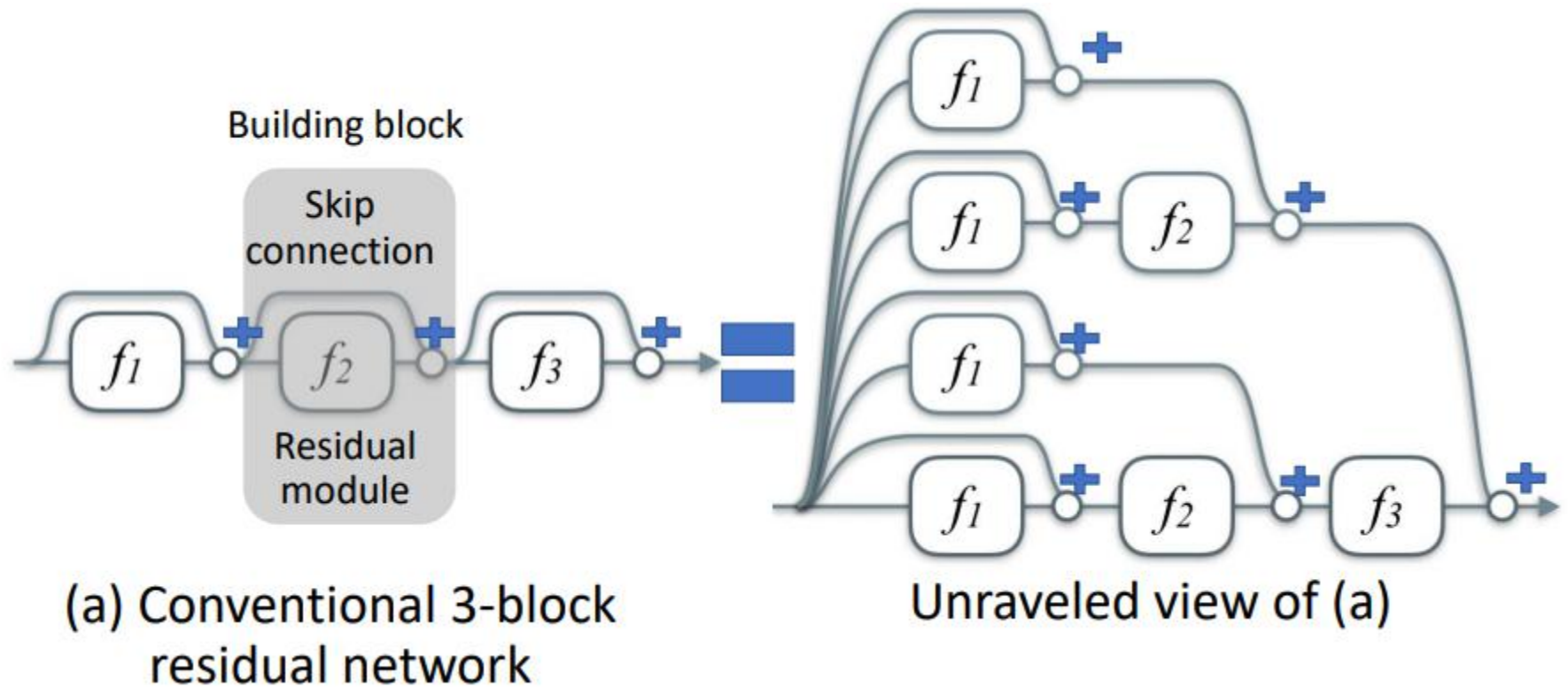
# Result on CIFAR-10



- Deep ResNets can be trained without difficulties
- Deeper ResNets have **lower training error**, and also lower test error

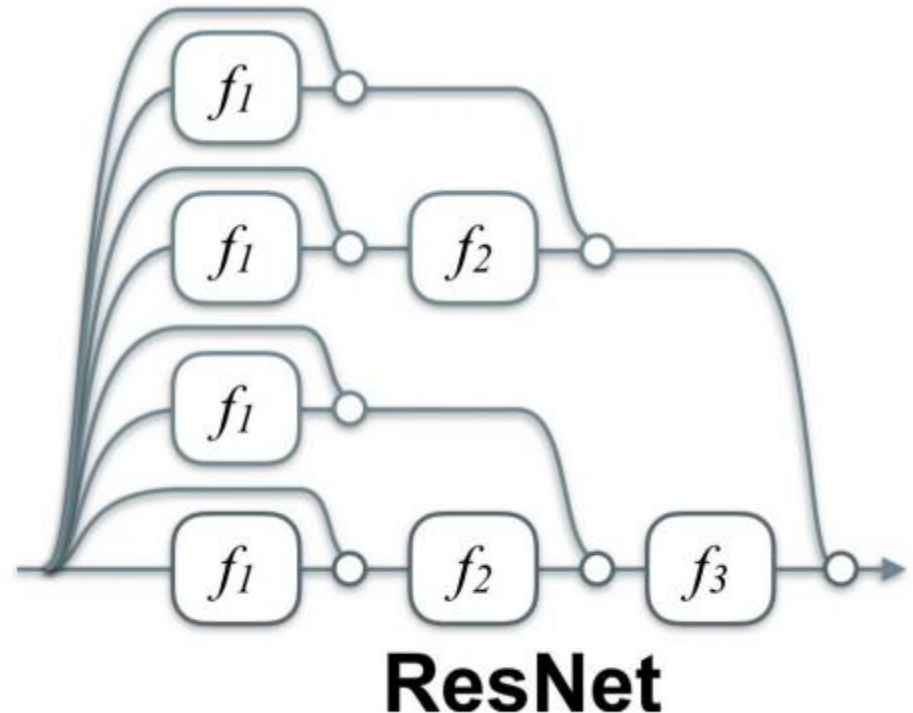# ImageNet experiments

- A practical design of going deeper



all-3x3  ⟷ similar complexity ⟷  bottleneck (for ResNet-50/101/152)

# Why does this work?



Building block
Skip connection
Residual module

$f_1$ $f_2$ $f_3$

(a) Conventional 3-block residual network

$f_1$
$f_1$ $f_2$
$f_1$
$f_1$ $f_2$ $f_3$

Unraveled view of (a)

# Why does this work?



The unraveled view is equivalent and showcases the many paths in ResNet.
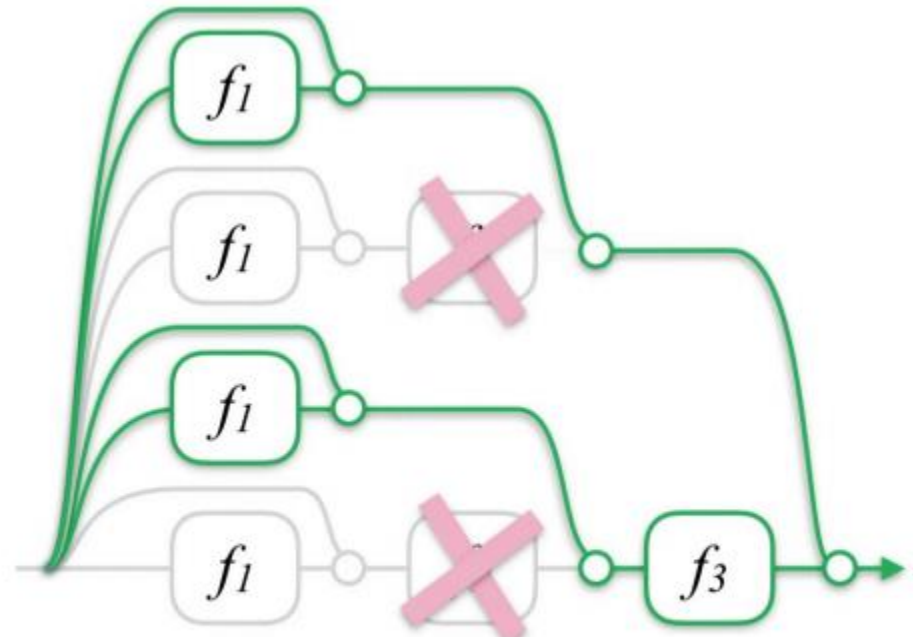
**VGG**

**ResNet**

# Deletion of one layer at test time
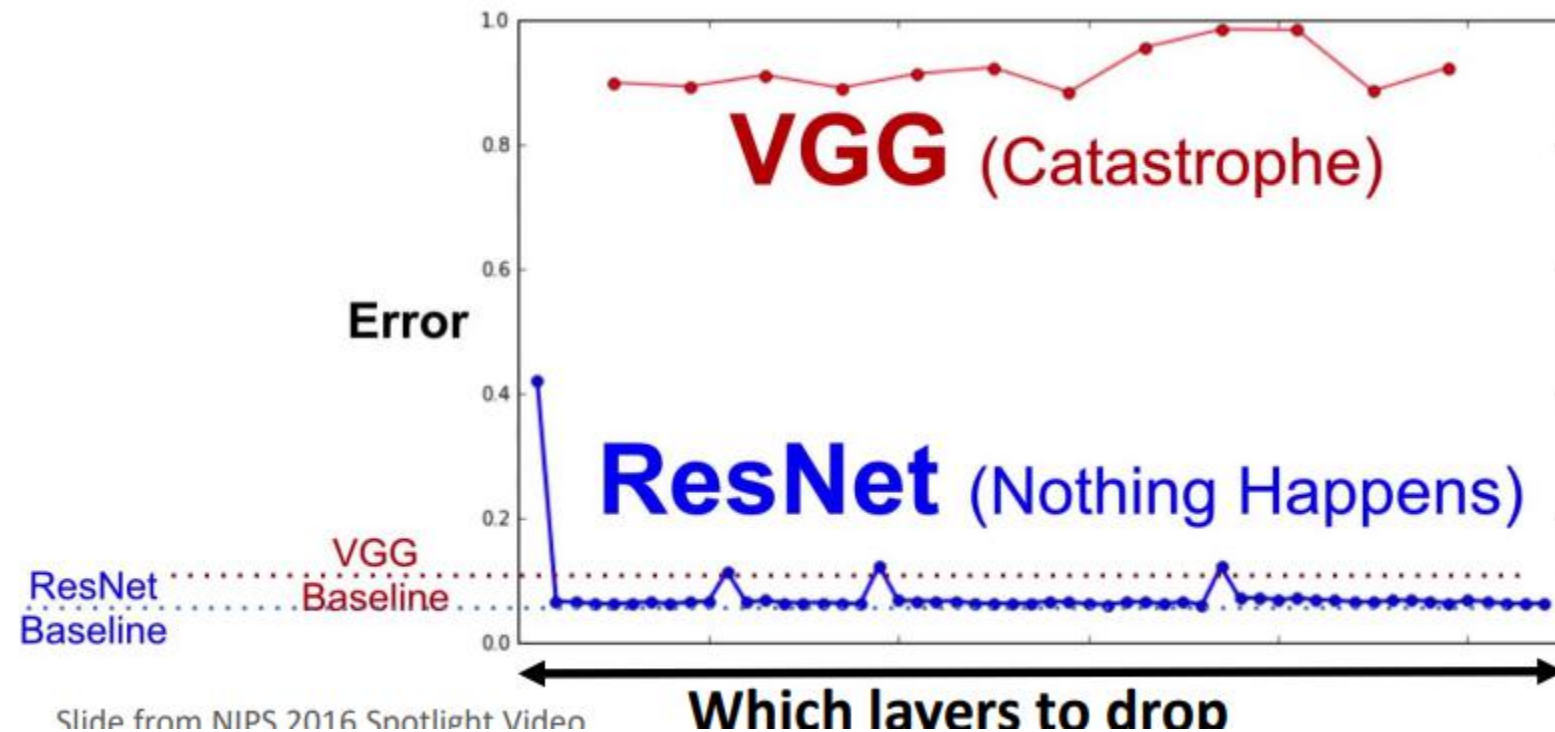


**VGG**
**All** paths are affected

**ResNet**
Only **half** of the paths are affected

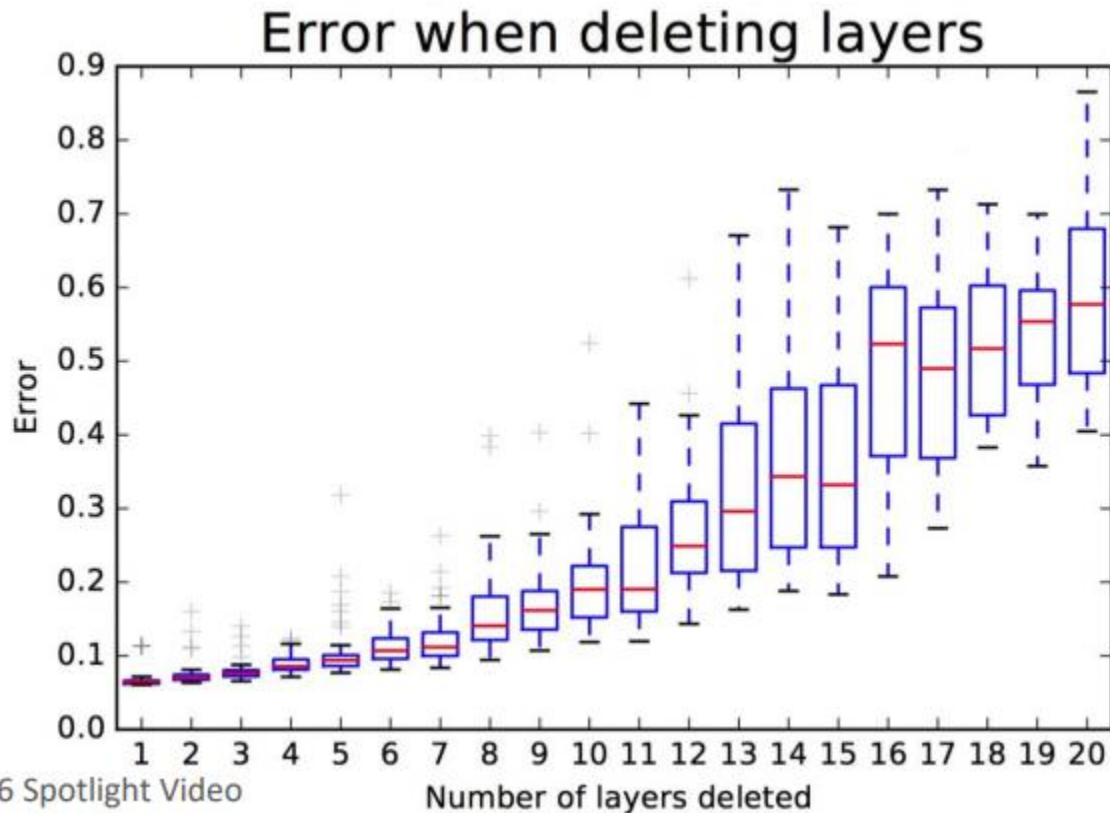# Deletion of one layer at test time



Slide from NIPS 2016 Spotlight Video

# Deletion of several layers



Error when deleting layers

# Conclusion 1

- Residual Networks consist of many paths.
- Although trained jointly, they do not strongly depend on each other: Ensemble-like behavior
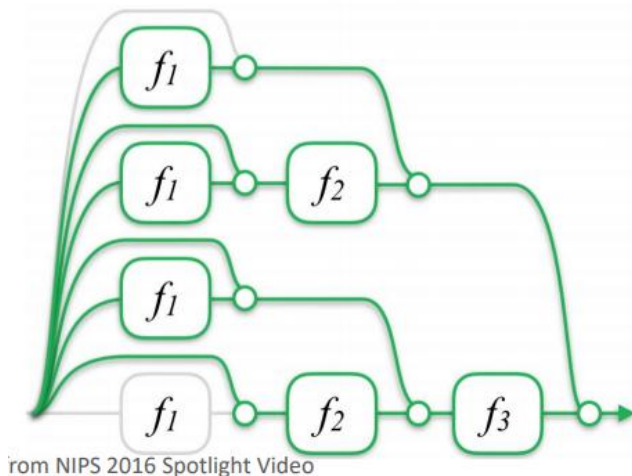
# Path Length

### Distribution of path length



There are very few **short paths...**

And very few **long paths...**

Most paths are **medium length!**
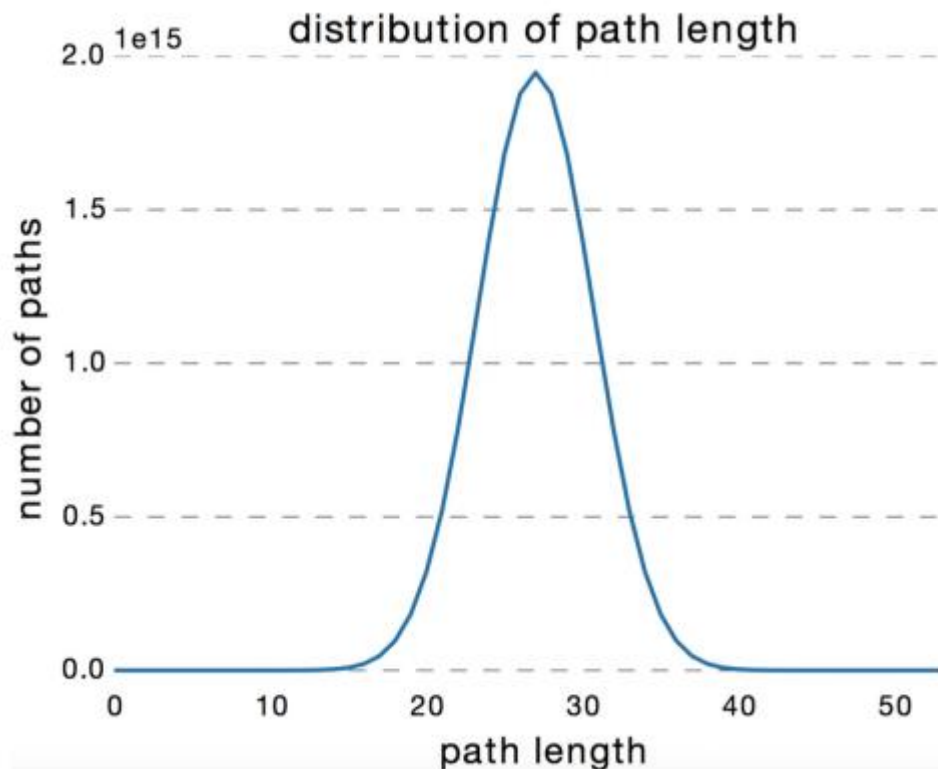
From NIPS 2016 Spotlight Video

**Residual networks contain many paths.**

Previous networks have a single path.

**Only short paths contribute gradient during training.**

Vanishing gradient suppresses gradient from long paths.

# Distribution of path lenth

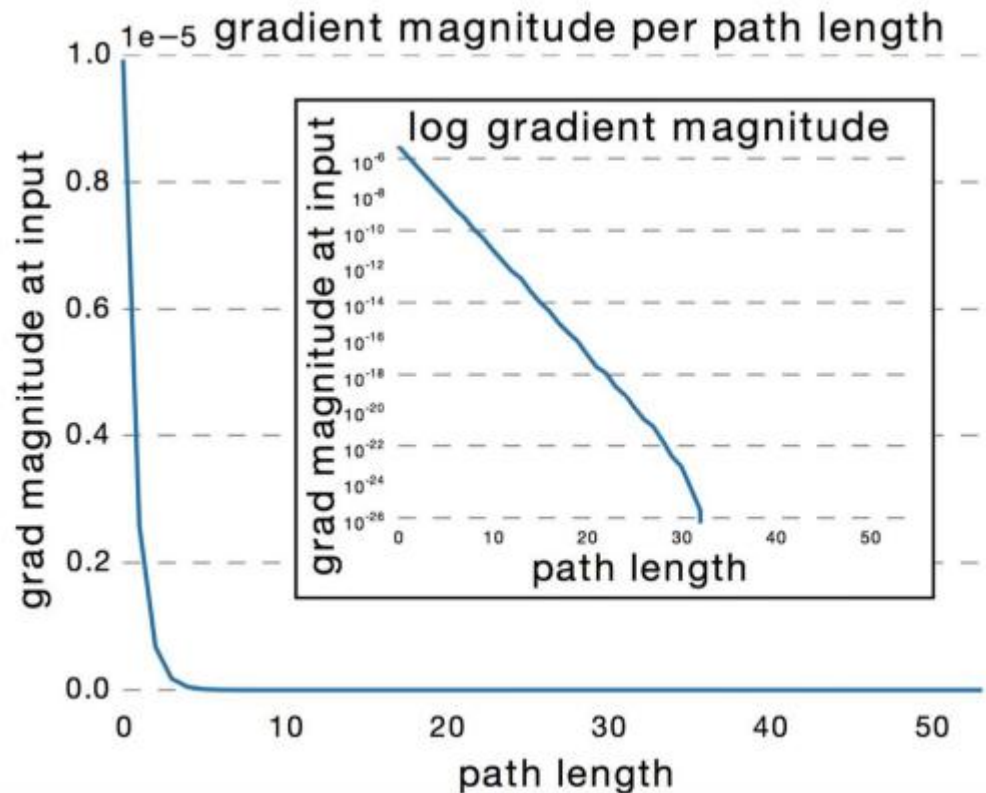

There are very few **short paths...**

And very few **long paths...**

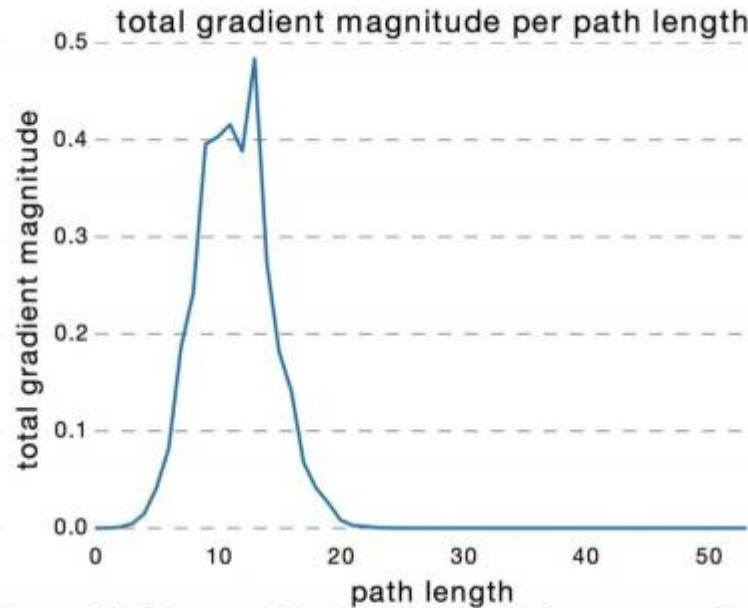Most paths are **medium length!**

Paths length follows a **binomial distribution.**

# Vanishing gradient

The gradient magnitude **decreases exponentially** with increasing path length.

# Gradient during training with path length



total gradient magnitude per path length

Combining the path length distribution and the vanishing gradients, one can observe that most of the gradient comes from relatively short paths.

# Conclusion 2

- Residual Networks consist of many paths.
- Although trained jointly, they do not strongly depend on each other: Ensemble-like behavior


- Most paths through a ResNet are relatively short.
- During training, gradients only flow through short paths.

Reference:

- Residual Networks Behave Like Ensembles of Relatively Shallow Networks

- Deep Residual Learning for Image Recognition