

# GoogLeNet

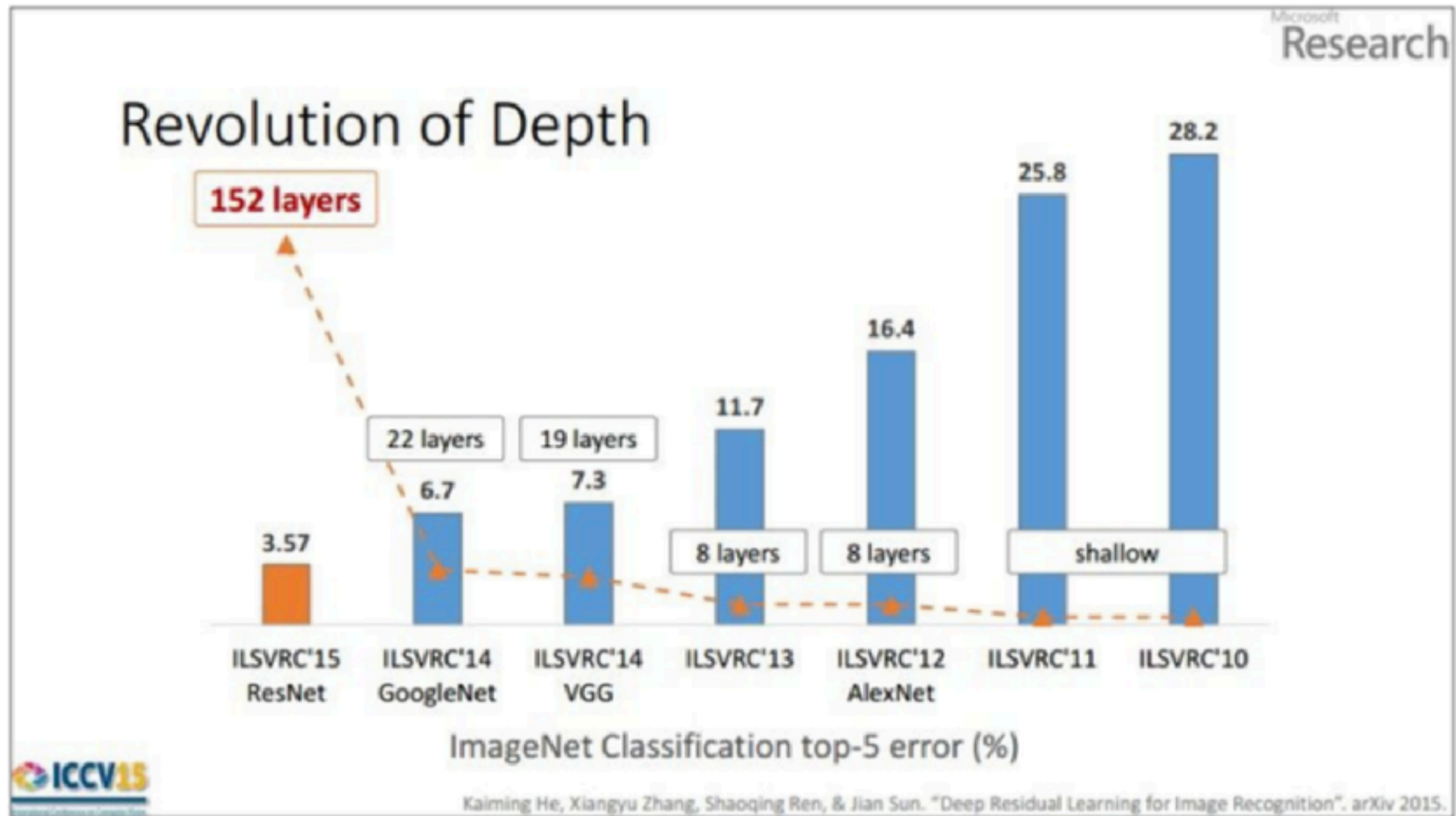
*Cho Sung Man*

# Index

- Introduction
- Architecture
- Results

# Introduction

# Introduction



A meme featuring Leonardo DiCaprio and Matt Damon from the movie Inception. DiCaprio is on the left, looking slightly to the right with a serious expression. Damon is on the right, leaning in towards DiCaprio. The background is a blurred office setting.

**WE NEED TO GO**

**DEEPER**

# Introduction

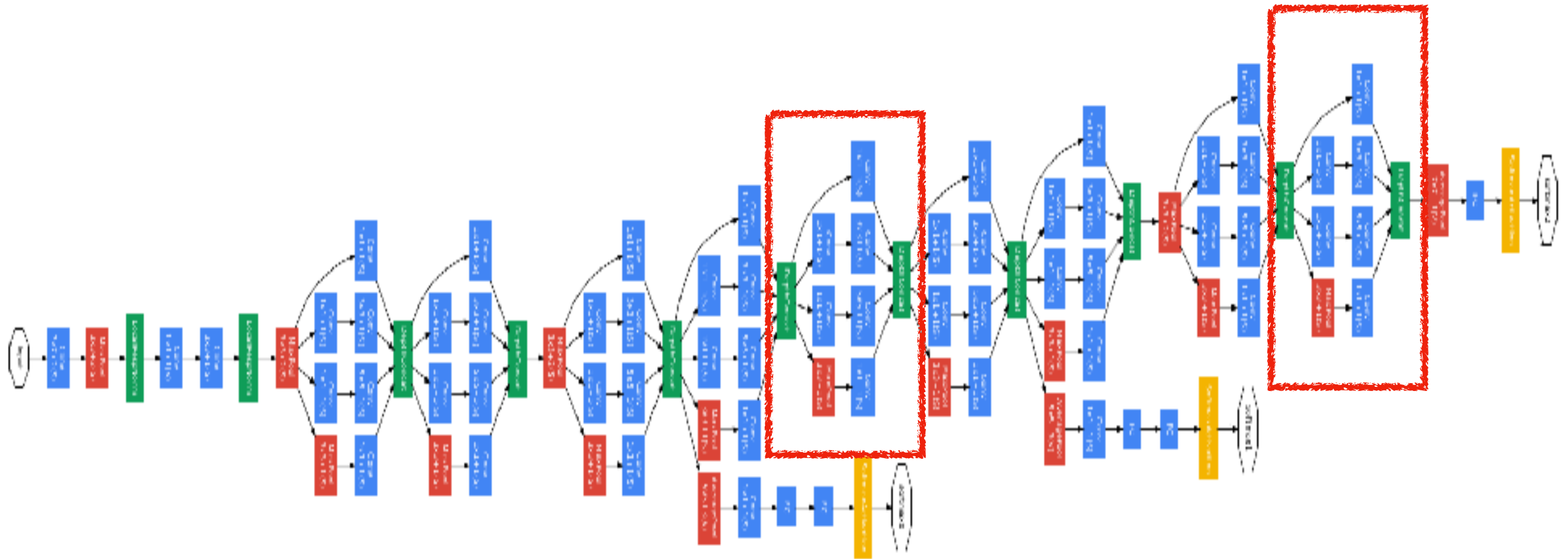


**Inception**

**Real Trade-Off Problem**

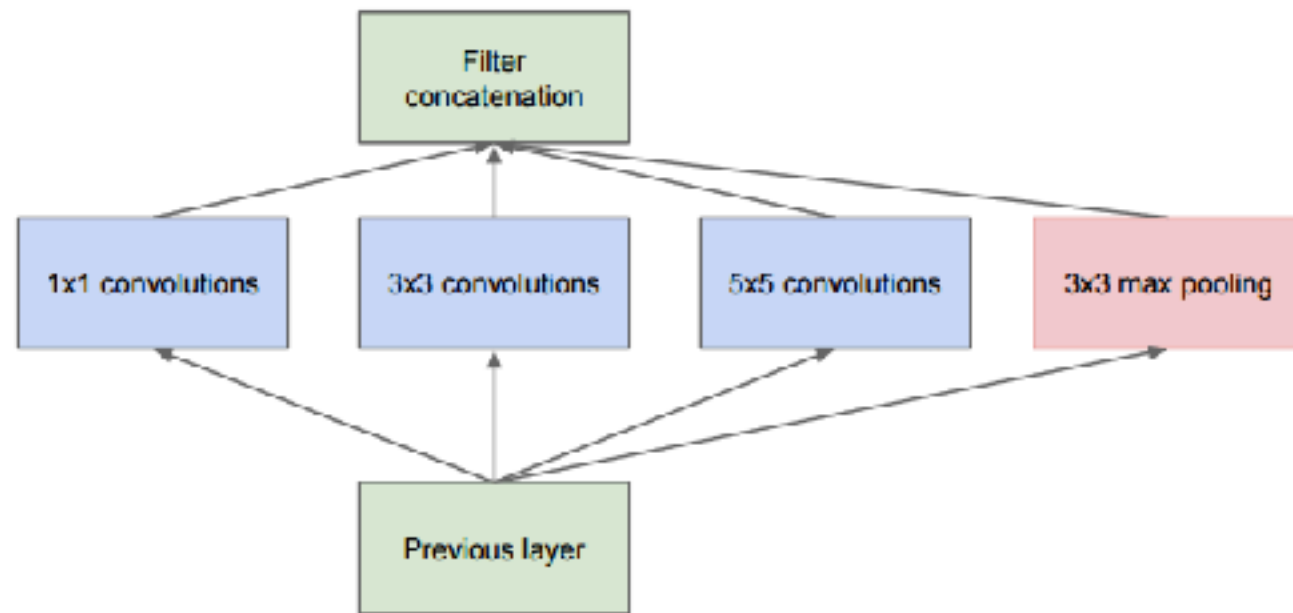
# Architecture

# Architecture

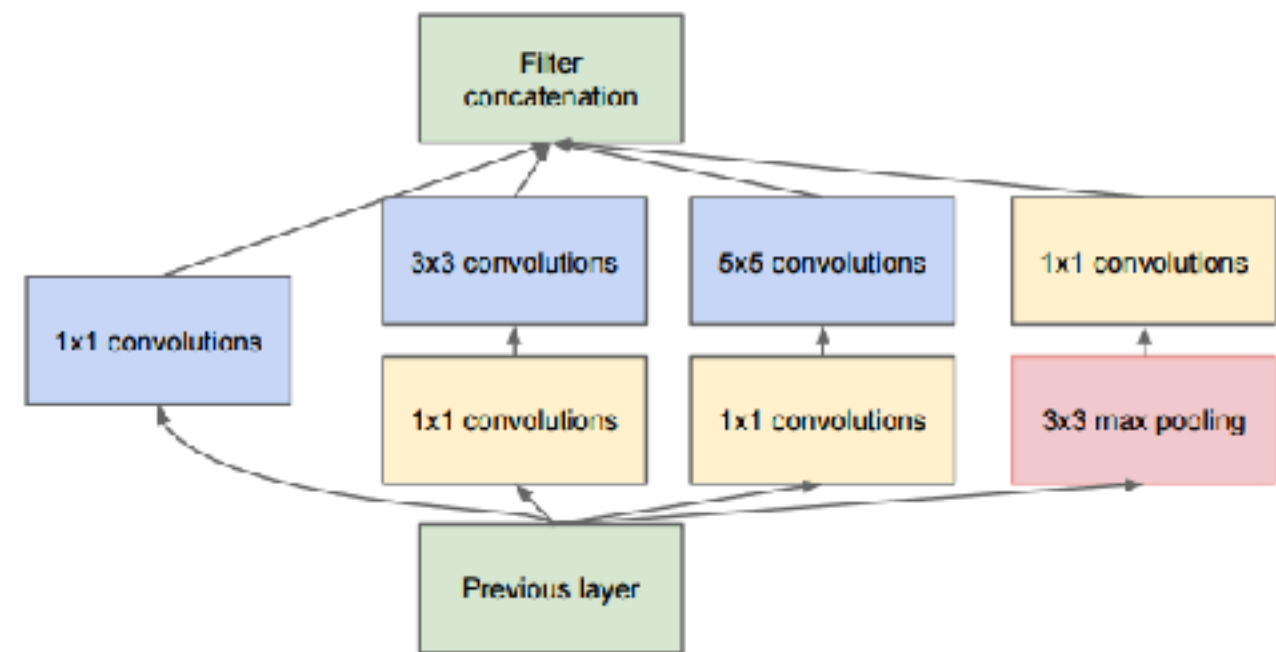




# Architecture



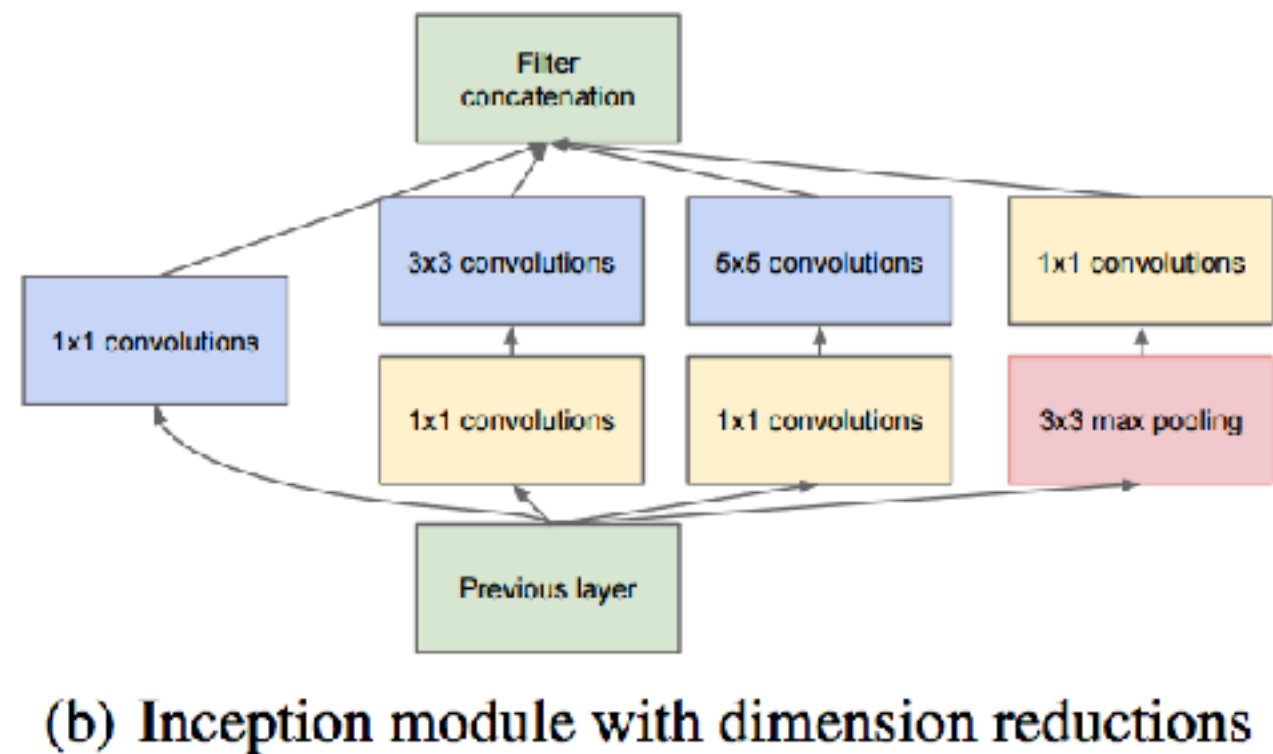
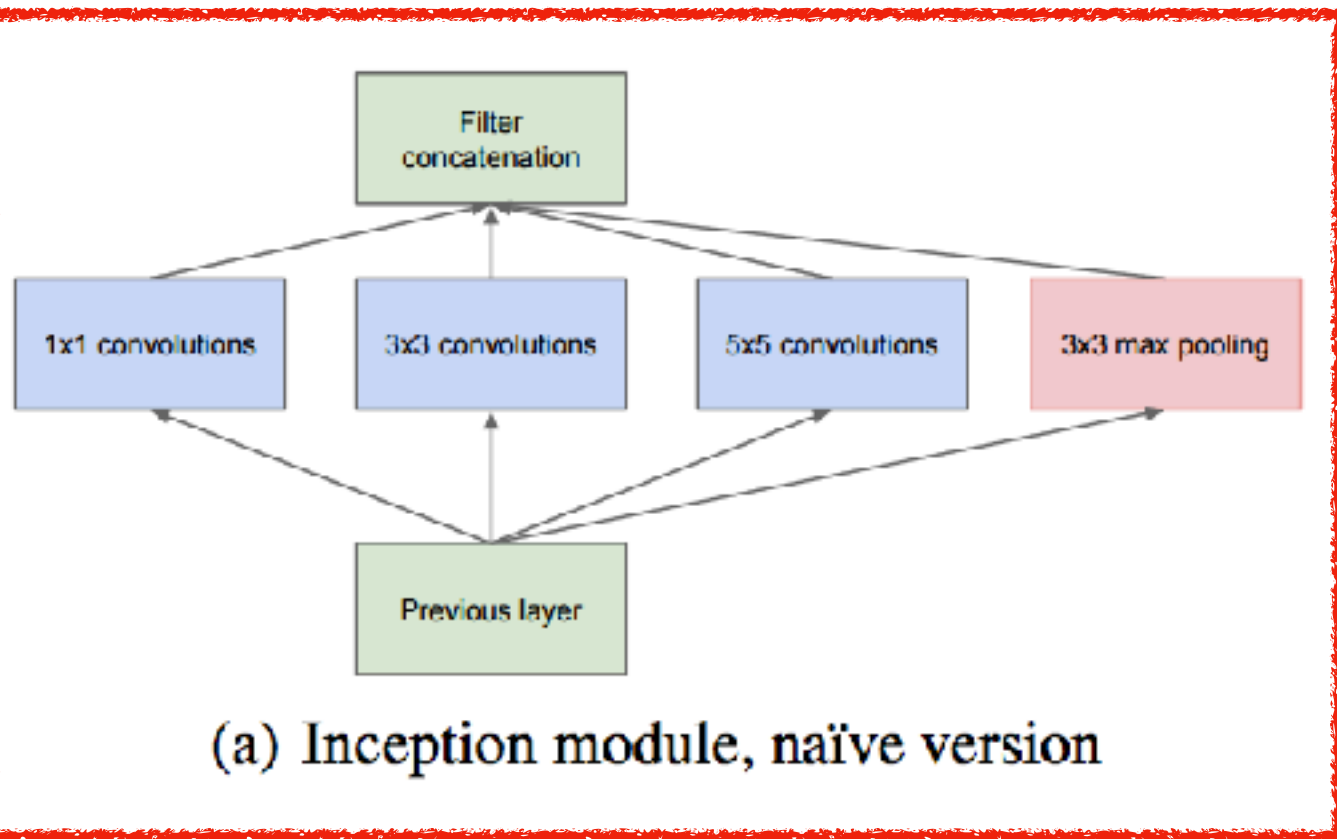
(a) Inception module, naïve version



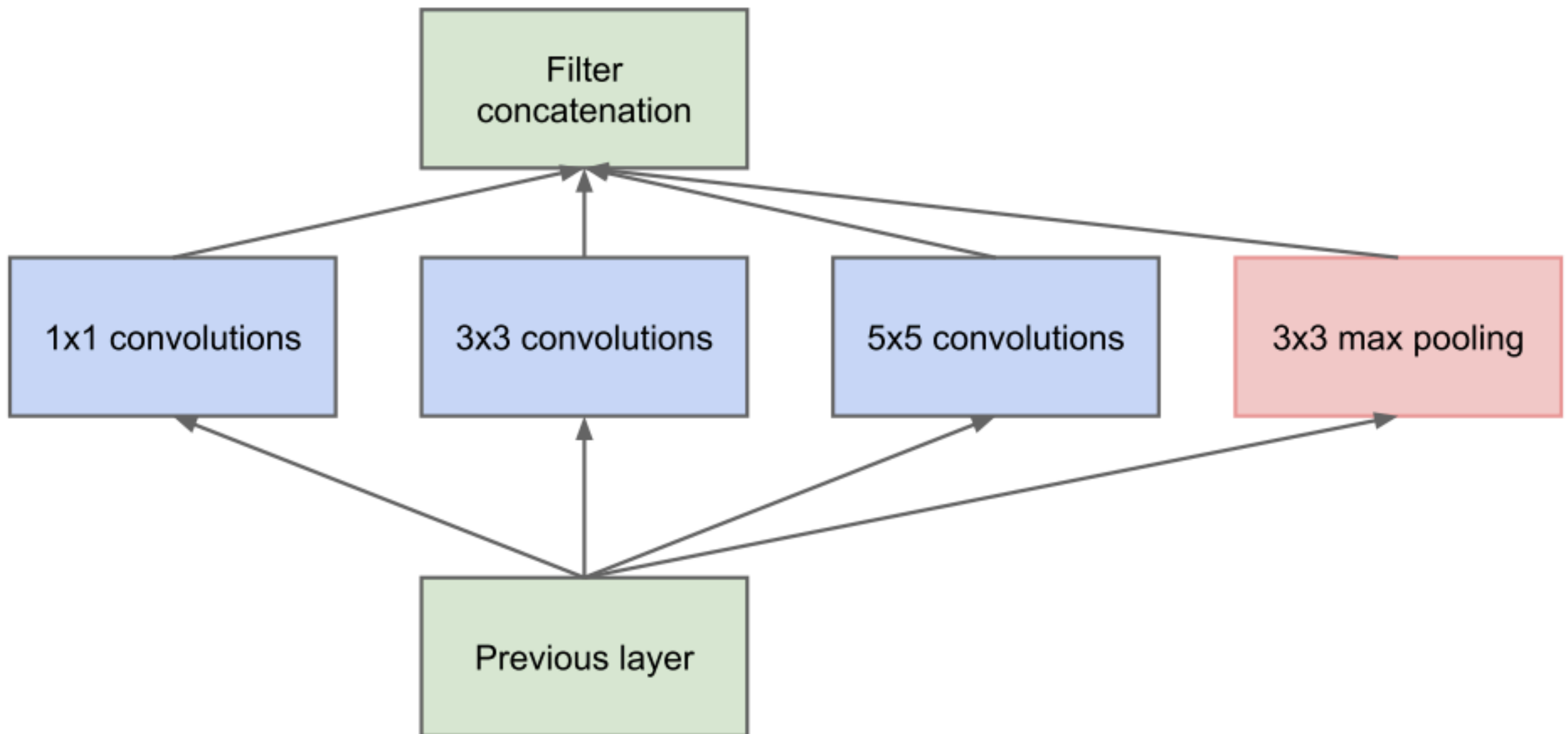
(b) Inception module with dimension reductions

Sparse Layer & Dense Matrix

# Architecture

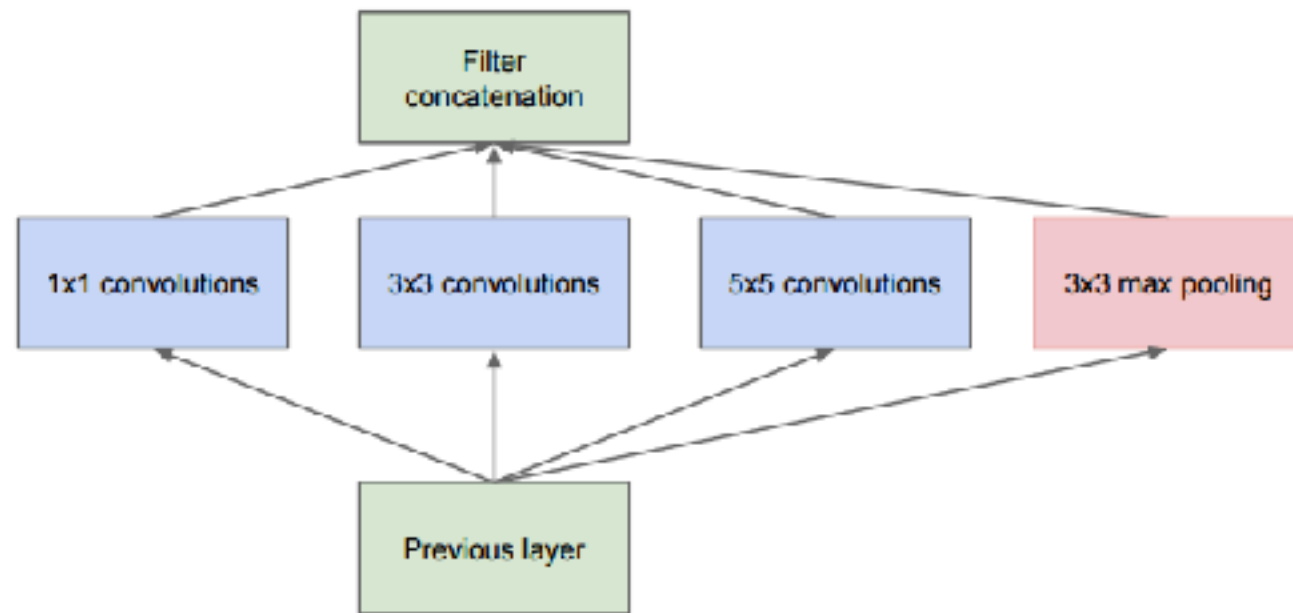


# Architecture

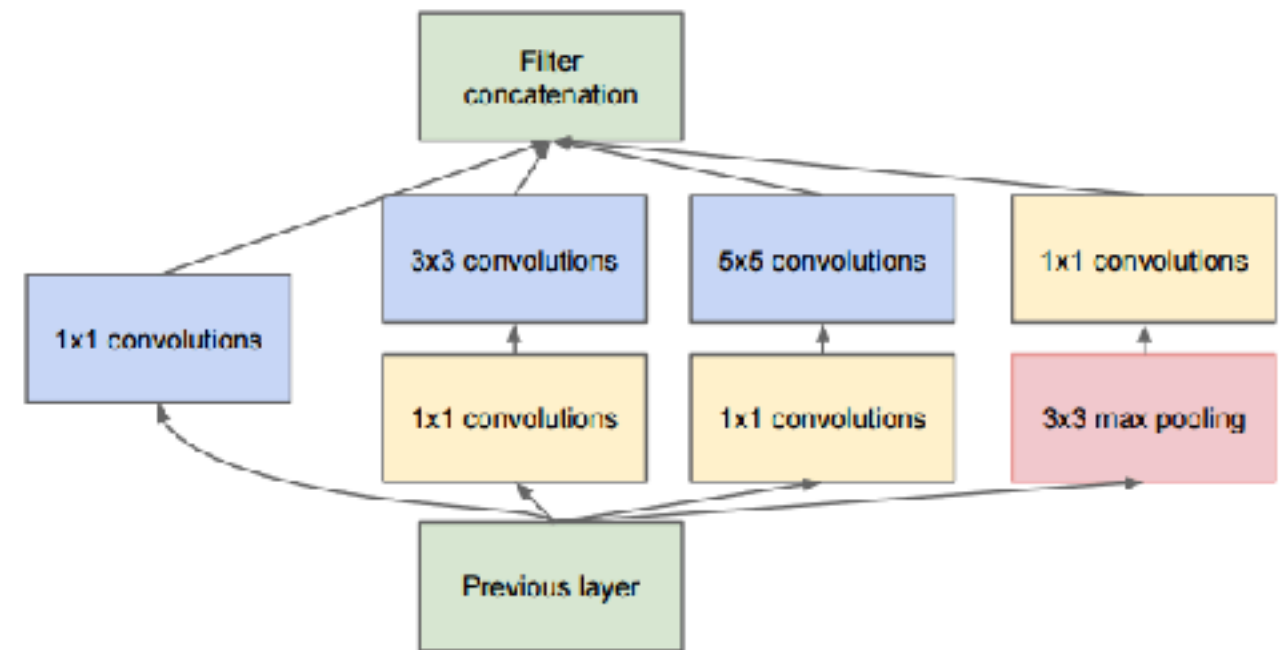


**Too Much Calculation.. So...**

# Architecture



(a) Inception module, naïve version



(b) Inception module with dimension reductions

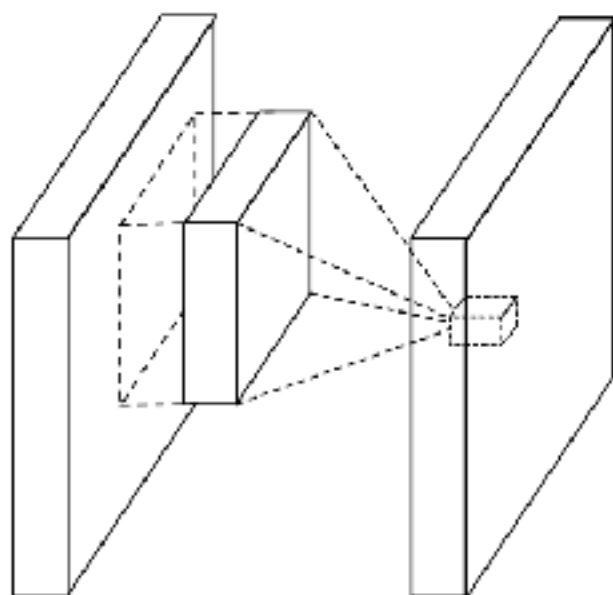
# Architecture

**NIN**

(Network In Network)

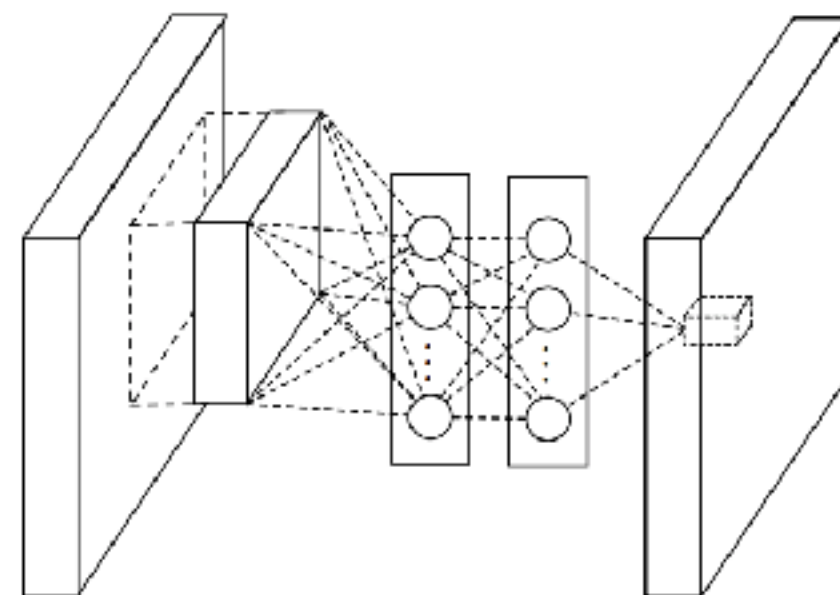
# Network In Network

**Conv -> Pooling**



(a) Linear convolution layer

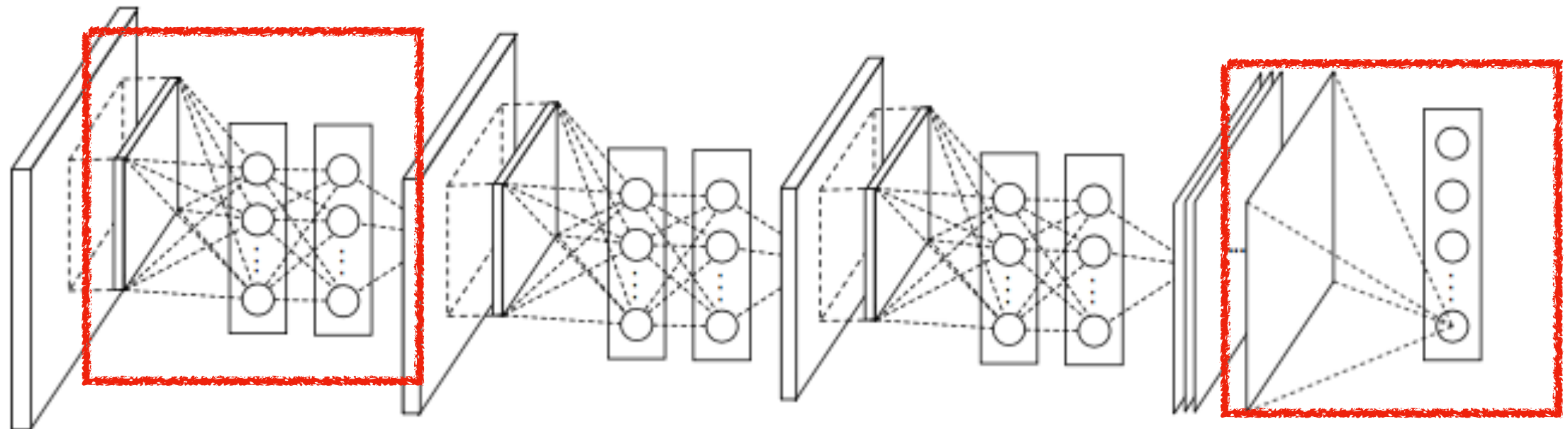
**Conv -> Mlp**



(b) Mlpconv layer

Figure 1: Comparison of linear convolution layer and mlpconv layer. The linear convolution layer includes a linear filter while the mlpconv layer includes a micro network (we choose the multilayer perceptron in this paper). Both layers map the local receptive field to a confidence value of the latent concept.

# Network In Network

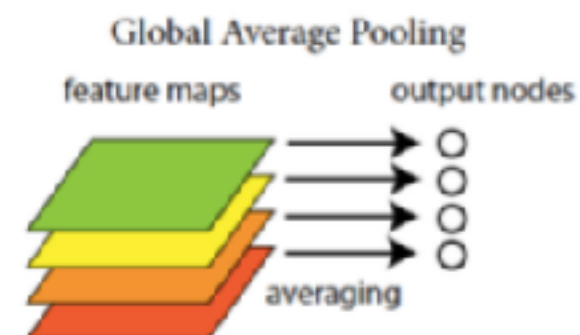
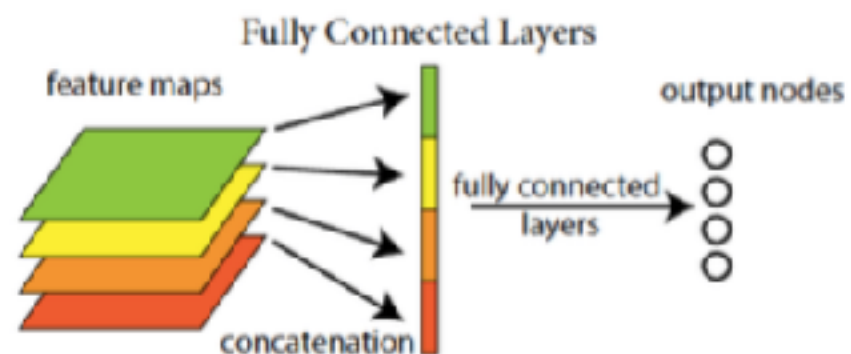


**Mlpconv Layer**

**Average Pooling**

**CNN**

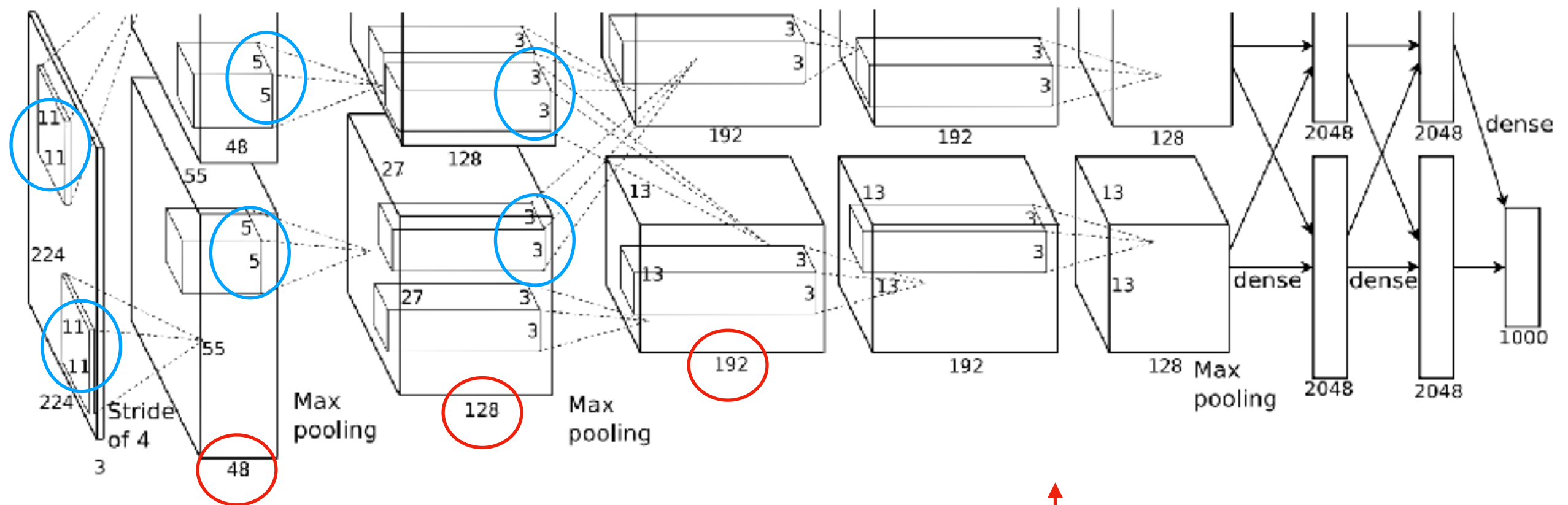
**NIN**



**Parameter Reduction**

# Network In Network

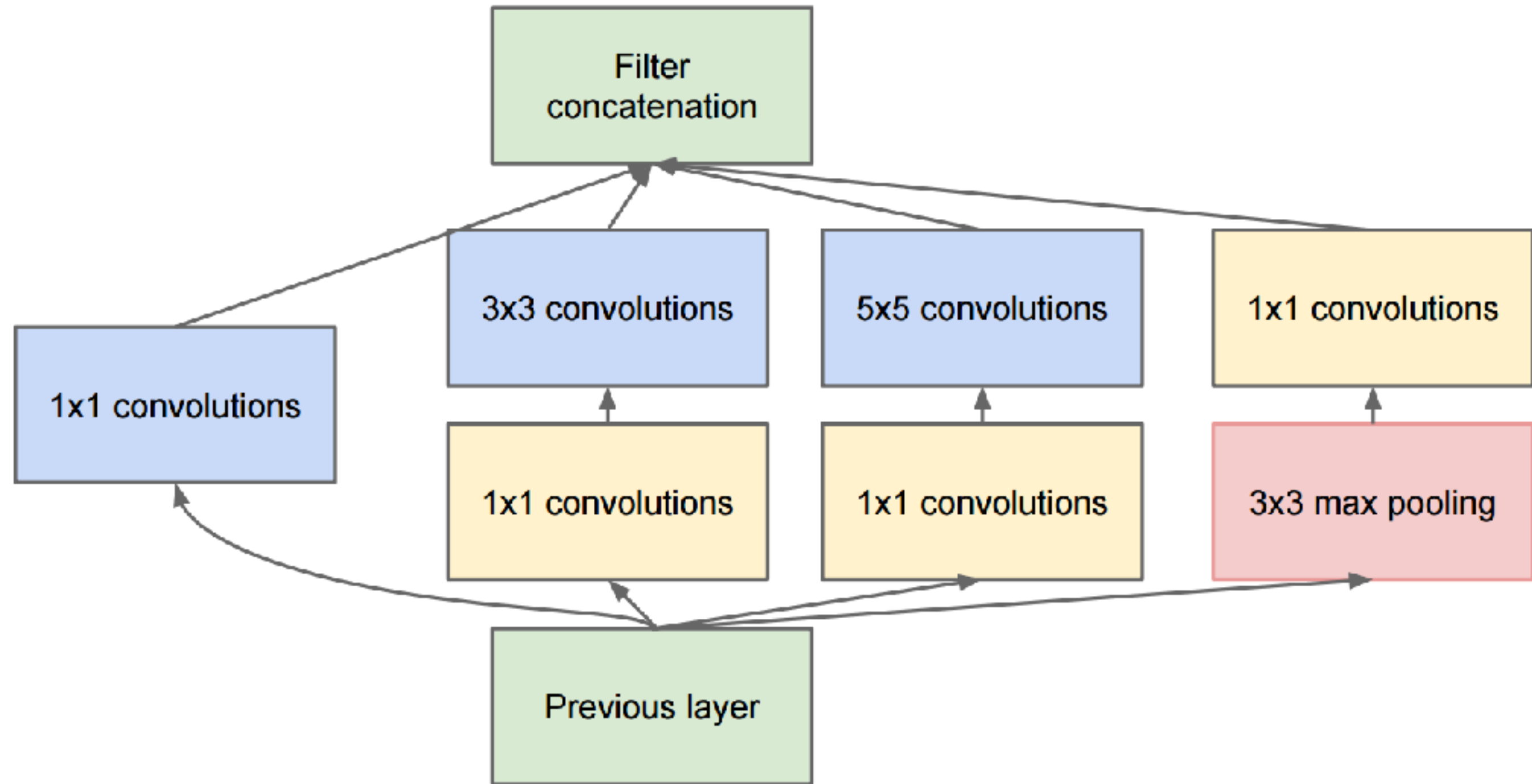
AlexNet



Convolution Layer : B, W, H, C

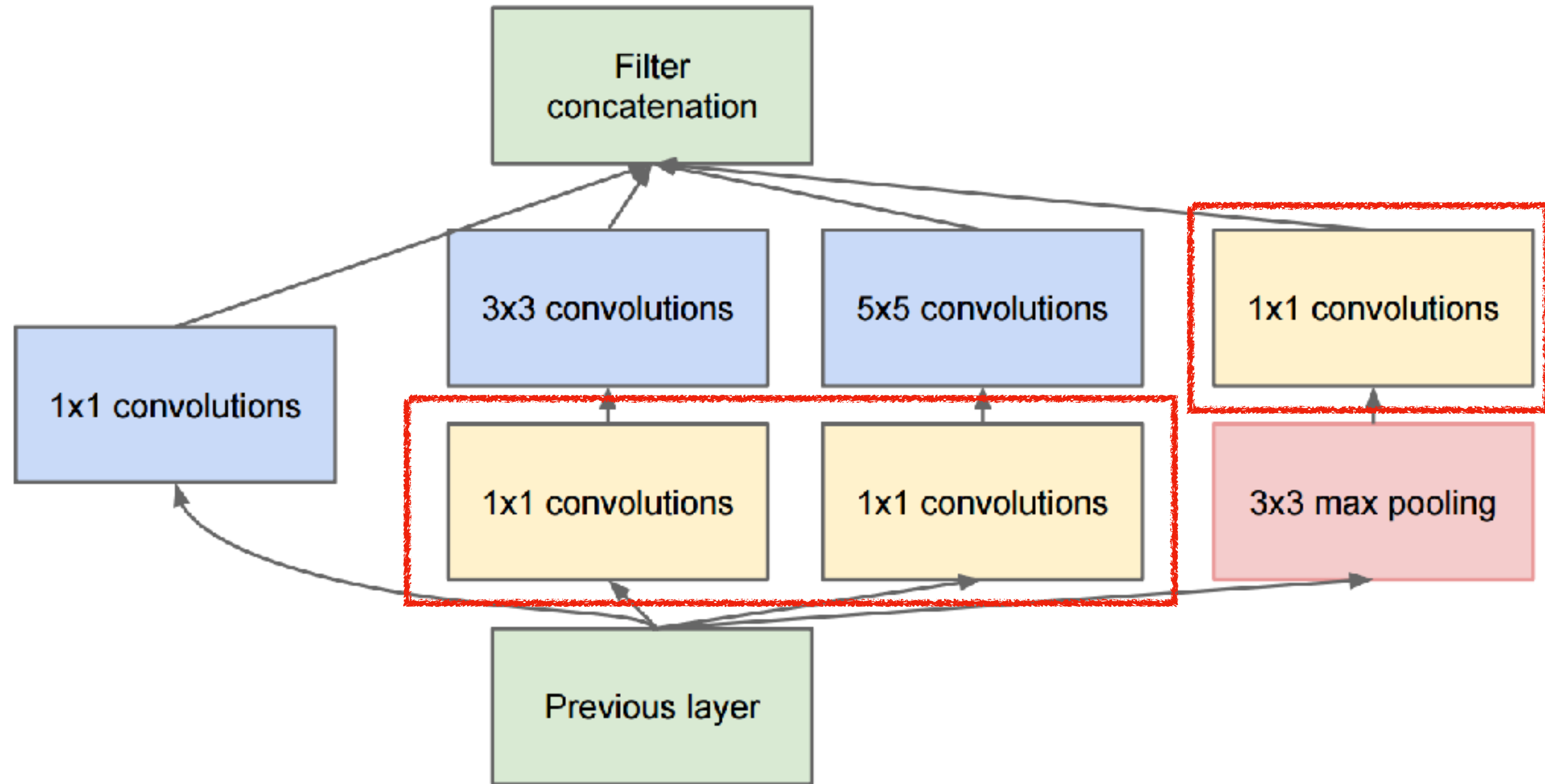


# Architecture



**Dimension Reductions**

# Architecture

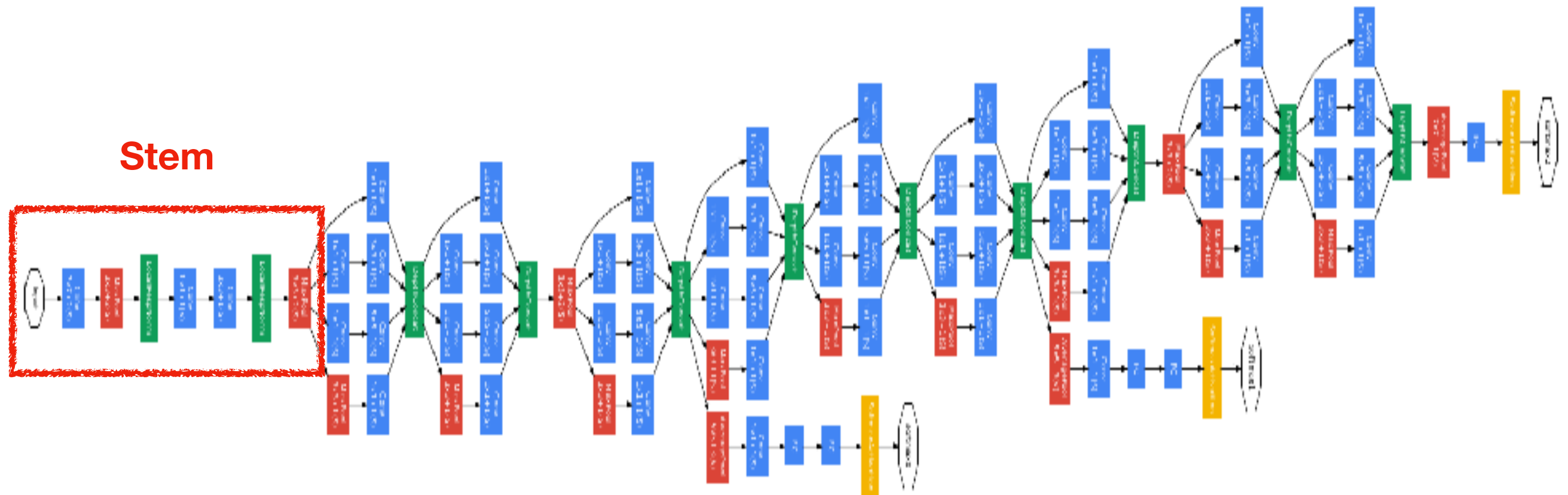


**Dimension Reductions**

**So Easy Concept !**  
**It's ALL !**

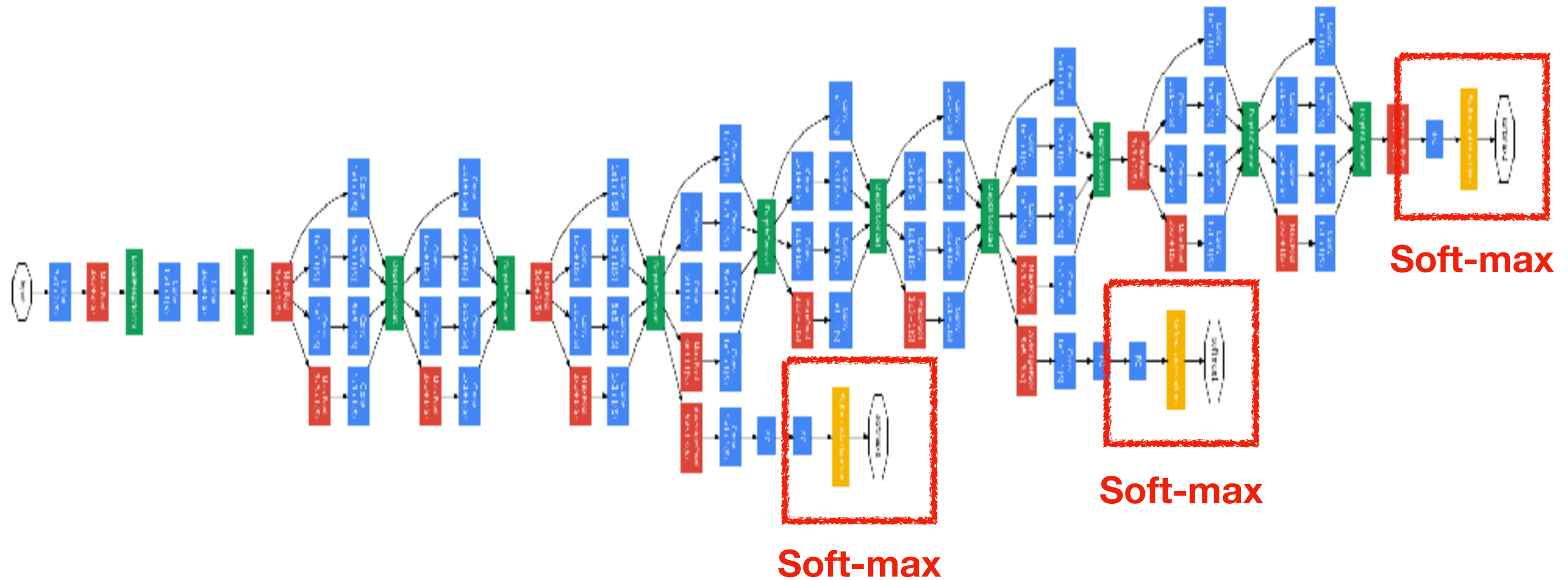
**NO!**

# Architecture



**In Early Layer, No Inception !**

# Architecture



**To prevent Vanishing, There are 3 Soft-max !**

# Results

# Results

type	patch size/ stride	output size	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	params	ops
convolution	7×7/2	112×112×64	1							2.7K	34M
max pool	3×3/2	56×56×64	0								
convolution	3×3/1	56×56×192	2		64	192				112K	360M
max pool	3×3/2	28×28×192	0								
inception (3a)		28×28×256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28×28×480	2	128	128	192	32	96	64	380K	304M
max pool	3×3/2	14×14×480	0								
inception (4a)		14×14×512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14×14×512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14×14×512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14×14×528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14×14×832	2	256	160	320	32	128	128	840K	170M
max pool	3×3/2	7×7×832	0								
inception (5a)		7×7×832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7×7×1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7×7/1	1×1×1024	0								
dropout (40%)		1×1×1024	0								
linear		1×1×1000	1							1000K	1M
softmax		1×1×1000	0								

22 Layers

Table 1: GoogLeNet incarnation of the Inception architecture.



# Results

Team	Year	Place	Error (top-5)	Uses external data
SuperVision	2012	1st	16.4%	no
SuperVision	2012	1st	15.3%	Imagenet 22k
Clarifai	2013	1st	11.7%	no
Clarifai	2013	1st	11.2%	Imagenet 22k
MSRA	2014	3rd	7.35%	no
VGG	2014	2nd	7.32%	no
GoogLeNet	2014	1st	6.67%	no

Table 2: Classification performance.

# Results

<b>Number of models</b>	<b>Number of Crops</b>	<b>Cost</b>	<b>Top-5 error</b>	<b>compared to base</b>
1	1	1	10.07%	base
1	10	10	9.15%	-0.92%
1	144	144	7.89%	-2.18%
7	1	7	8.09%	-1.98%
7	10	70	7.62%	-2.45%
7	144	1008	6.67%	-3.45%

Table 3: GoogLeNet classification performance break down.

# Results

<b>Team</b>	<b>Year</b>	<b>Place</b>	<b>mAP</b>	<b>external data</b>	<b>ensemble</b>	<b>approach</b>
UvA-Eurovision	2013	1st	22.6%	none	?	Fisher vectors
Deep Insight	2014	3rd	40.5%	ImageNet 1k	3	CNN
CUHK DeepID-Net	2014	2nd	40.7%	ImageNet 1k	?	CNN
GoogLeNet	2014	1st	43.9%	ImageNet 1k	6	CNN

Table 4: Comparison of detection performances. Unreported values are noted with question marks.

# Results

<b>Team</b>	<b>mAP</b>	<b>Contextual model</b>	<b>Bounding box regression</b>
Trimps-Soushen	31.6%	no	?
Berkeley Vision	34.5%	no	yes
UvA-Eurovision	35.4%	?	?
CUHK DeepID-Net2	37.7%	no	?
GoogLeNet	38.02%	no	no
Deep Insight	40.2%	yes	yes

Table 5: Single model performance for detection.

**Thank You !**