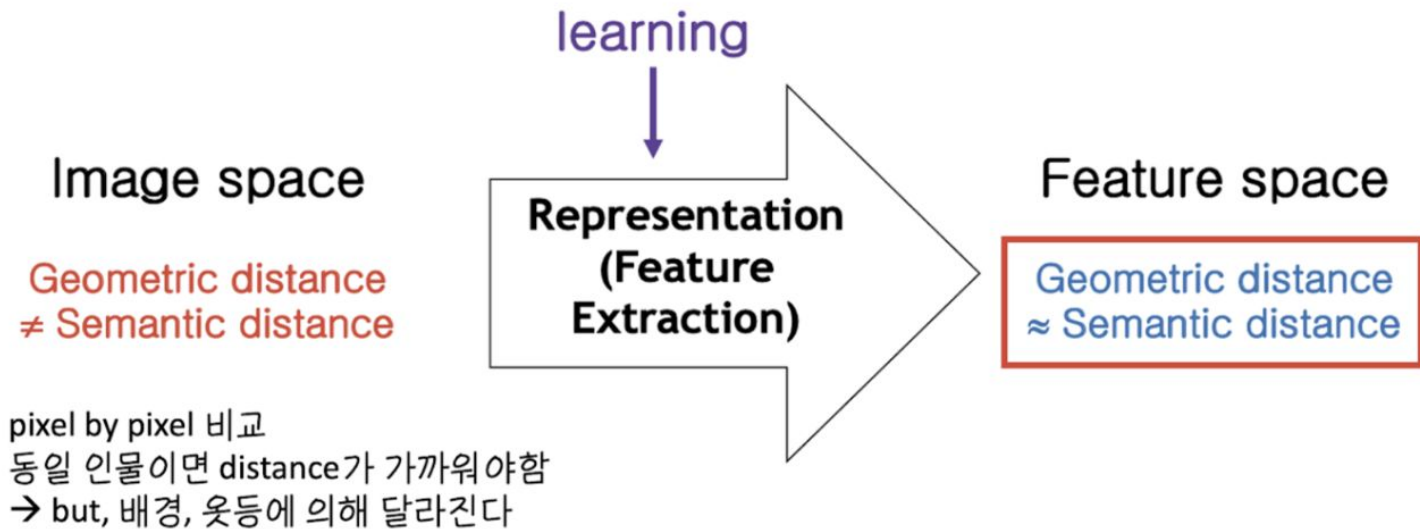


FaceNet: A Unified Embedding for Face Recognition and Clustering

(Florian Schroff, Dmitry Kalenichenko, James Philbin, 2015)

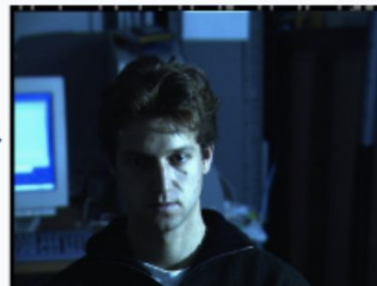
Background - Metric Learning



Abstract & Introduction



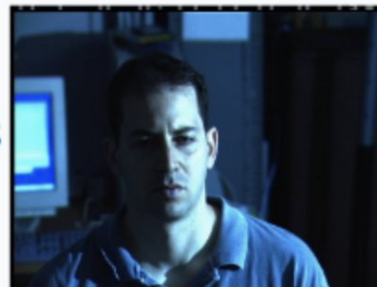
1.22



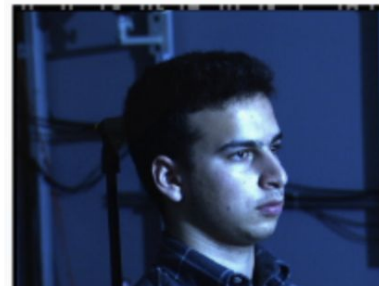
1.33



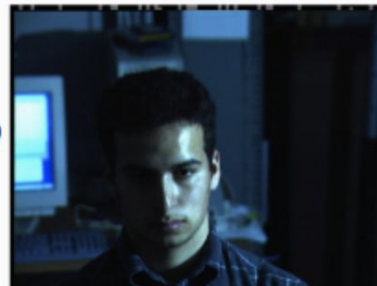
1.33



1.26

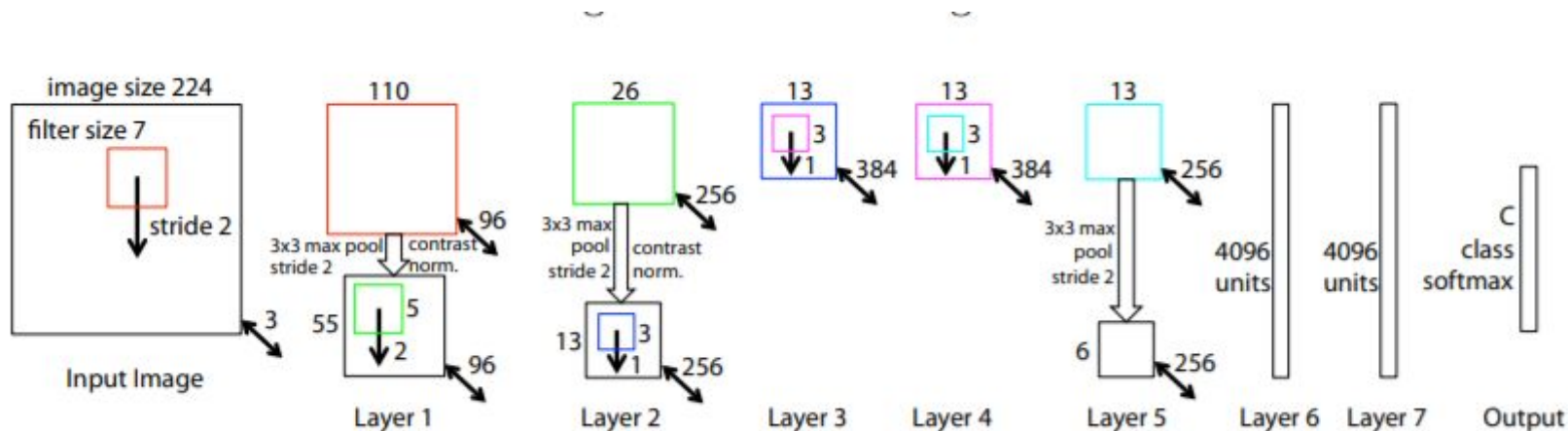


0.99



Related Work - based model

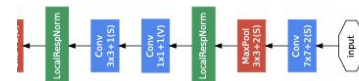
- Visualizing and understanding convolutional networks



Visualizing and Understanding Convolutional Networks, Matthew D. Zeiler and Rob Fergus
Dept. of Computer Science, New York University, USA, 2014

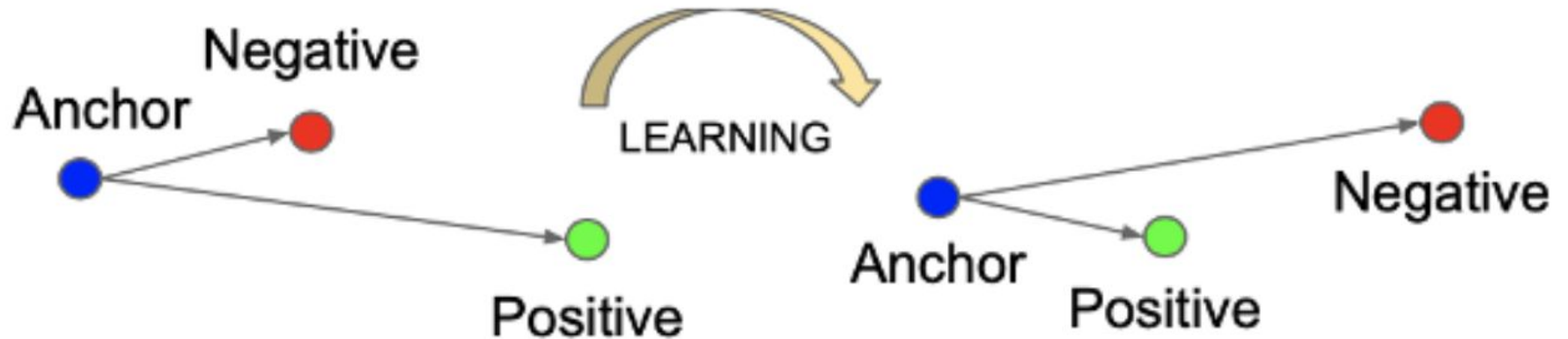
Related Work - based model - Inception model

- Going deeper with convolutions



Going deeper with convolutions, Szegedy, Dept of Computer Vision and Pattern Recognition, Cornell University, 2014

Method - Triplet Loss



Method - Triplet Loss

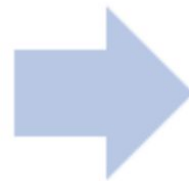
$$\|f(x_i^a) - f(x_i^p)\|_2^2 + \alpha < \|f(x_i^a) - f(x_i^n)\|_2^2$$

$$\text{Loss} = \sum_i^N \left[\|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right]$$

Method - Triplet Selection

$$\operatorname{argmax}_{x_i^p} \|f(x_i^a) - f(x_i^p)\|_2^2 \quad \text{Hard positive}$$

$$\operatorname{argmin}_{x_i^n} \|f(x_i^a) - f(x_i^n)\|_2^2 \quad \text{Hard negative}$$



계산량 ↑
오버피팅 문제

Method - Deep Convolutional Networks - NN1

layer	size-in	size-out	kernel	param	FLPS
conv1	220×220×3	110×110×64	7×7×3, 2	9K	115M
pool1	110×110×64	55×55×64	3×3×64, 2	0	
rnorm1	55×55×64	55×55×64		0	
conv2a	55×55×64	55×55×64	1×1×64, 1	4K	13M
conv2	55×55×64	55×55×192	3×3×64, 1	111K	335M
rnorm2	55×55×192	55×55×192		0	
pool2	55×55×192	28×28×192	3×3×192, 2	0	
conv3a	28×28×192	28×28×192	1×1×192, 1	37K	29M
conv3	28×28×192	28×28×384	3×3×192, 1	664K	521M
pool3	28×28×384	14×14×384	3×3×384, 2	0	
conv4a	14×14×384	14×14×384	1×1×384, 1	148K	29M
conv4	14×14×384	14×14×256	3×3×384, 1	885K	173M
conv5a	14×14×256	14×14×256	1×1×256, 1	66K	13M
conv5	14×14×256	14×14×256	3×3×256, 1	590K	116M
conv6a	14×14×256	14×14×256	1×1×256, 1	66K	13M
conv6	14×14×256	14×14×256	3×3×256, 1	590K	116M
pool4	14×14×256	7×7×256	3×3×256, 2	0	
concat	7×7×256	7×7×256		0	
fc1	7×7×256	1×32×128	maxout p=2	103M	103M
fc2	1×32×128	1×32×128	maxout p=2	34M	34M
fc7128	1×32×128	1×1×128		524K	0.5M
L2	1×1×128	1×1×128		0	
total				140M	1.6B

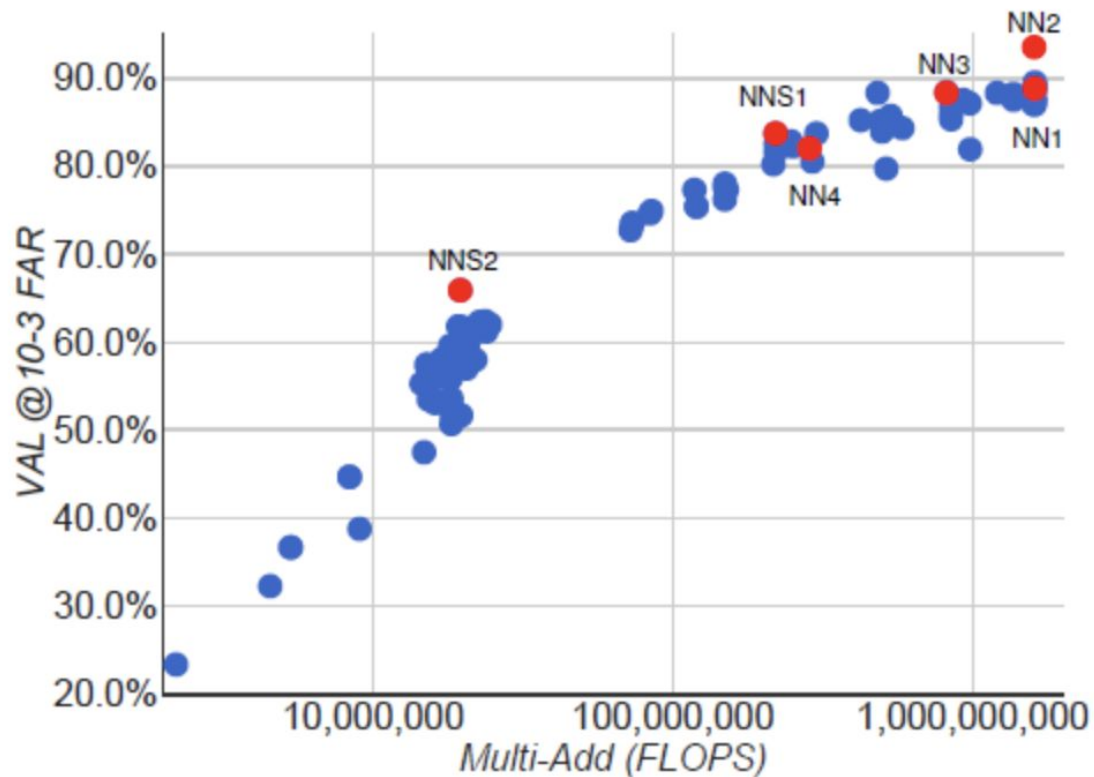
Method - Deep Convolutional Networks - NN2

[illegible]

Method - Deep Convolutional Networks

architecture	VAL
NN1 (Zeiler&Fergus 220×220)	$87.9\% \pm 1.9$
NN2 (Inception 224×224)	$89.4\% \pm 1.6$
NN3 (Inception 160×160)	$88.3\% \pm 1.7$
NN4 (Inception 96×96)	$82.0\% \pm 2.3$
NNS1 (mini Inception 165×165)	$82.4\% \pm 2.4$
NNS2 (tiny Inception 140×116)	$51.9\% \pm 2.9$

Method - Deep Convolutional Networks



Datasets and Evaluation

$$\text{TA}(d) = \{(i, j) \in \mathcal{P}_{\text{same}} \mid D(x_i, x_j) \leq d\}$$

$$\text{FA}(d) = \{(i, j) \in \mathcal{P}_{\text{diff}} \mid D(x_i, x_j) \leq d\}$$

$$\text{VAL}(d) = \frac{|\text{TA}(d)|}{|\mathcal{P}_{\text{same}}|}, \text{FAR}(d) = \frac{|\text{FA}(d)|}{|\mathcal{P}_{\text{diff}}|}$$

Datasets and Evaluation

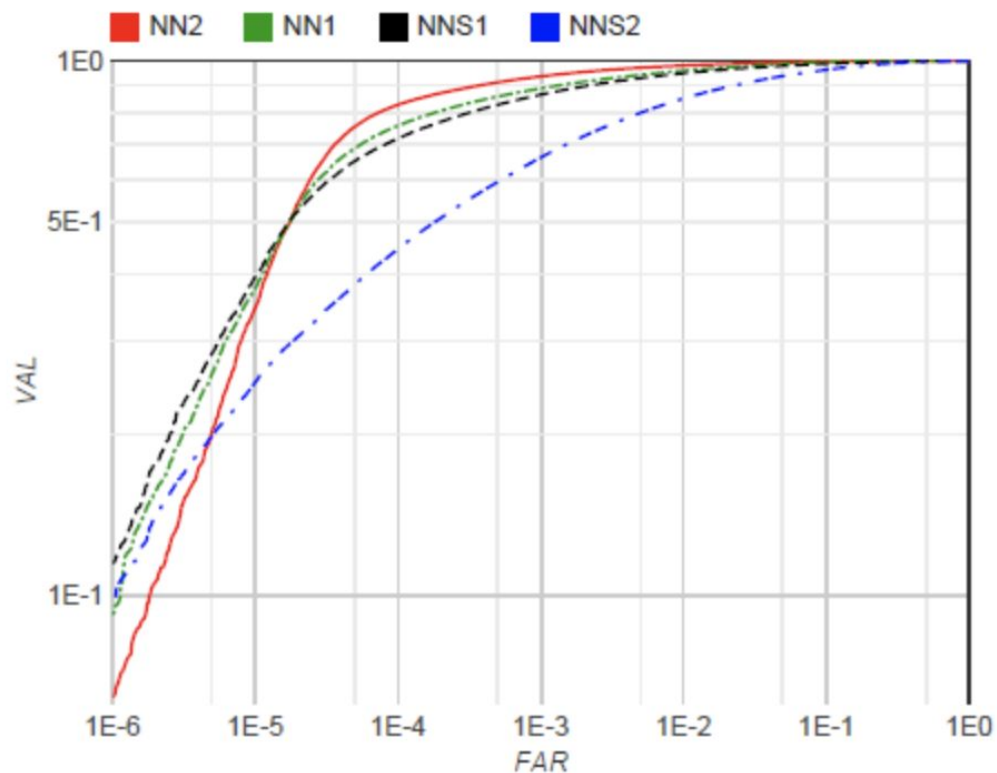
- **Hold-out Data Set**
- **Personal Photos**
- **Academic Datasets**

Experiments - Computation Accuracy Trade-off

정확도 , FLOPS : Trade-off 관계



Experiments - Effect of CNN Model



Experiments - Sensitivity to Image Quality

jpeg q	val-rate
10	67.3%
20	81.4%
30	83.9%
50	85.5%
70	86.1%
90	86.5%

#pixels	val-rate
1,600	37.8%
6,400	79.5%
14,400	84.5%
25,600	85.7%
65,536	86.4%

Experiments - Embedding Dimensionality

#dims	VAL
64	86.8% \pm 1.7
128	87.9% \pm 1.9
256	87.7% \pm 1.9
512	85.6% \pm 2.0

Experiments - Amount of Training Data

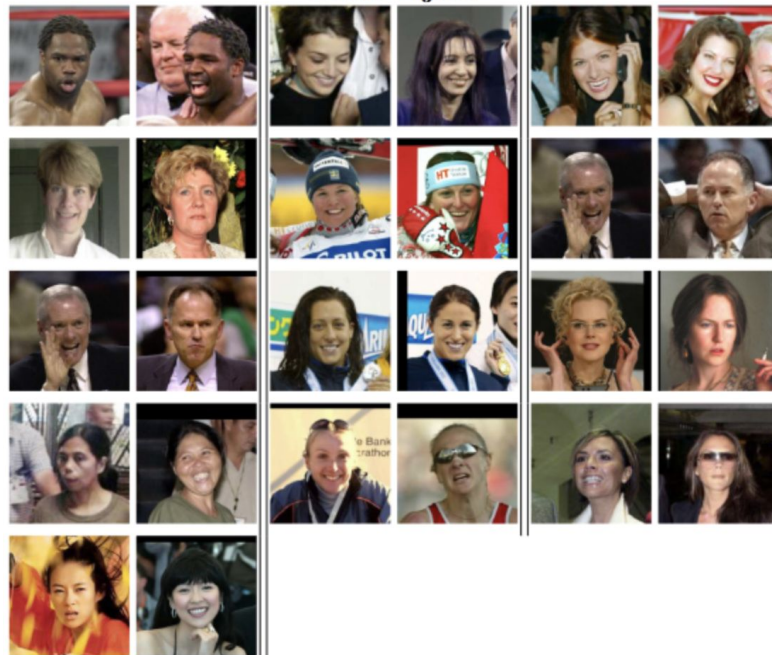
#training images	VAL
2,600,000	76.3%
26,000,000	85.1%
52,000,000	85.1%
260,000,000	86.2%

Experiments - Performance on LFW

False accept



False reject



Experiments - Face Clustering

