

facebook

Web-Scale Training for Face Identification

Yaniv Taigman
Facebook AI Research

FACE RECOGNITION

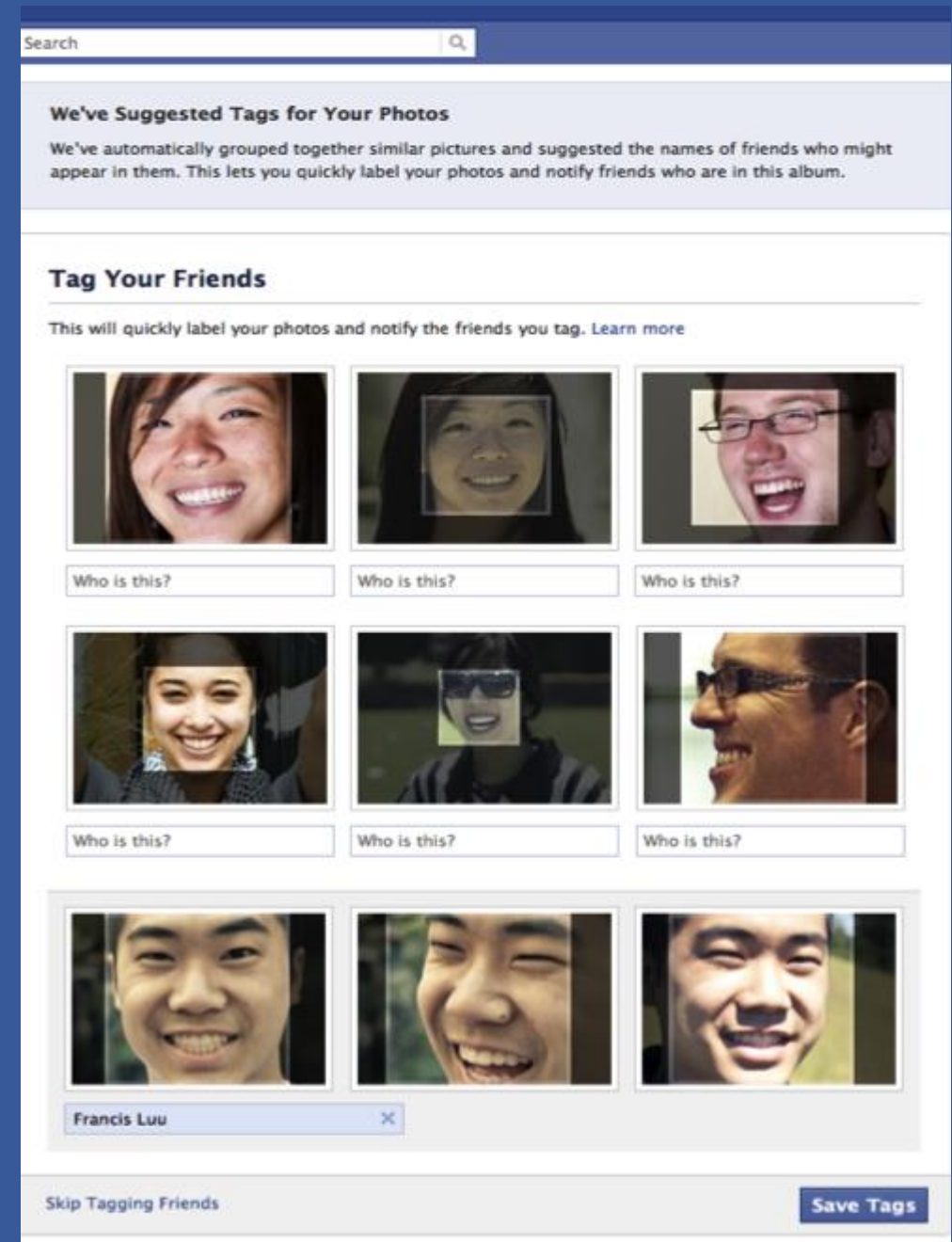


Why faces?



1. One class. Billions of **unique** instances.
2. Plays an important role in our social interactions, conveying people's identity; **The most frequent entity** in the media by far: e.g. ~1.2 faces / Photo by avg
1. Enables many **applications** in Man-Machine interaction

Applications



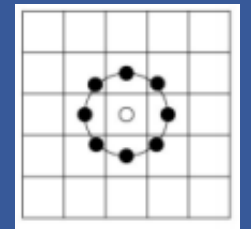
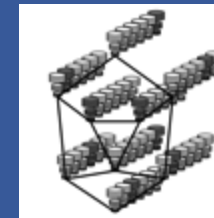
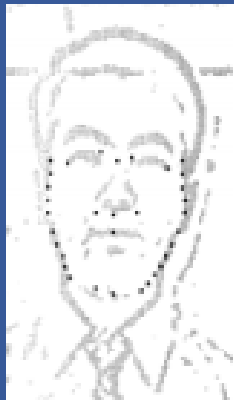
Face Recognition main objective

Find a representation & similarity measure such that:

- Intra-subject similarity is high
- Inter-subject similarity is low



Milestones in Face Recognition



1964

**BLED SOE
FACE
RECOGNITION**

1973

**KANADE'S
THESIS**

1991

**TURK &
PENTLAND
EIGENFACES**

1997

**BELHUMEUR
FISHERFACE**

1999

**BLANZ &
VETTER
MORPHABLE
FACES**

1999

**WISKOTT
EBGM**

2001

**VIOLA &
JONES
BOOSTING**

2006

**AHONEN
LBP**

Problem solved?

NIST's best-performer's on:

1. Its internal dataset with 1.6 million identities: 95.9%
2. On LFW (public) with 'only' 4,249 identities: 56.7%

→ Answer: No.

- L. Best-Rowden, H. Han, C. Otto, B. Klare, and A. K. Jain. Unconstrained face recognition: Identifying a person of interest from a media collection. TR MSU-CSE-14-1, 2014.

Types of Face Recognition

- ‘Constrained’ – Mainly for traditional purposes
- ‘Unconstrained’ – General purpose

CONSTRAINED



NIST'S FR VENDOR TEST (FRVT) 2006

UNCONSTRAINED



IN THE WILD

Challenges in Unconstrained Face Recognition

1. POSE

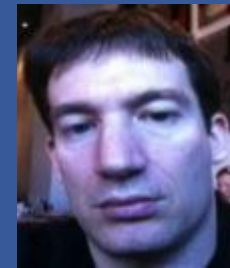
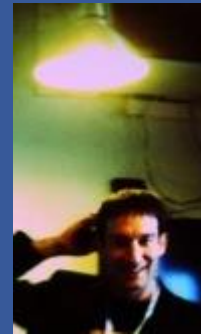
2. ILLUMINATION

3. EXPRESSION

4. AGING

5. OCCLUSION

Probes for example



Gallery



Unconstrained Face Recognition Era: The Labeled Faces in the Wild (LFW)



13,233 PHOTOS OF 5,749



Labeled faces in the wild: A database for studying face recognition in unconstrained environments, Huang, Jain, Learned-Miller, ECCVW, 2008

LFW: Progress over the recent 7 years

- Labeled faces in the wild: A database for studying face recognition in unconstrained environments, ECCVW, 2008.
- Descriptor methods in the Wild, ECCV-W 2008
- Attribute and simile classifiers for face verification, ICCV 2009.
- Multiple one-shots for utilizing class label information, BMVC 2009.
- Large scale strongly supervised ensemble metric learning, with applications to face verification and retrieval, NEC Labs TR, 2012.
- Learning hierarchical representations for face verification with convolutional deep belief networks, CVPR, 2012.
- Bayesian face revisited: A joint formulation, ECCV 2012.
- Tom-vs-pete classifiers and identity preserving alignment for face verification, BMVC 2012.
- Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification, CVPR 2013.
- Probabilistic elastic matching for pose variant face verification, CVPR 2013.
- Fusing robust face region descriptors via multiple metric learning for face recognition in the wild, CVPR 2013.
- Fisher vector faces in the wild, BMVC 2013.
- Hybrid deep learning for computing face similarities, ICCV 2013.
- A practical transfer learning algorithm for face verification, ICCV 2013.

Verification



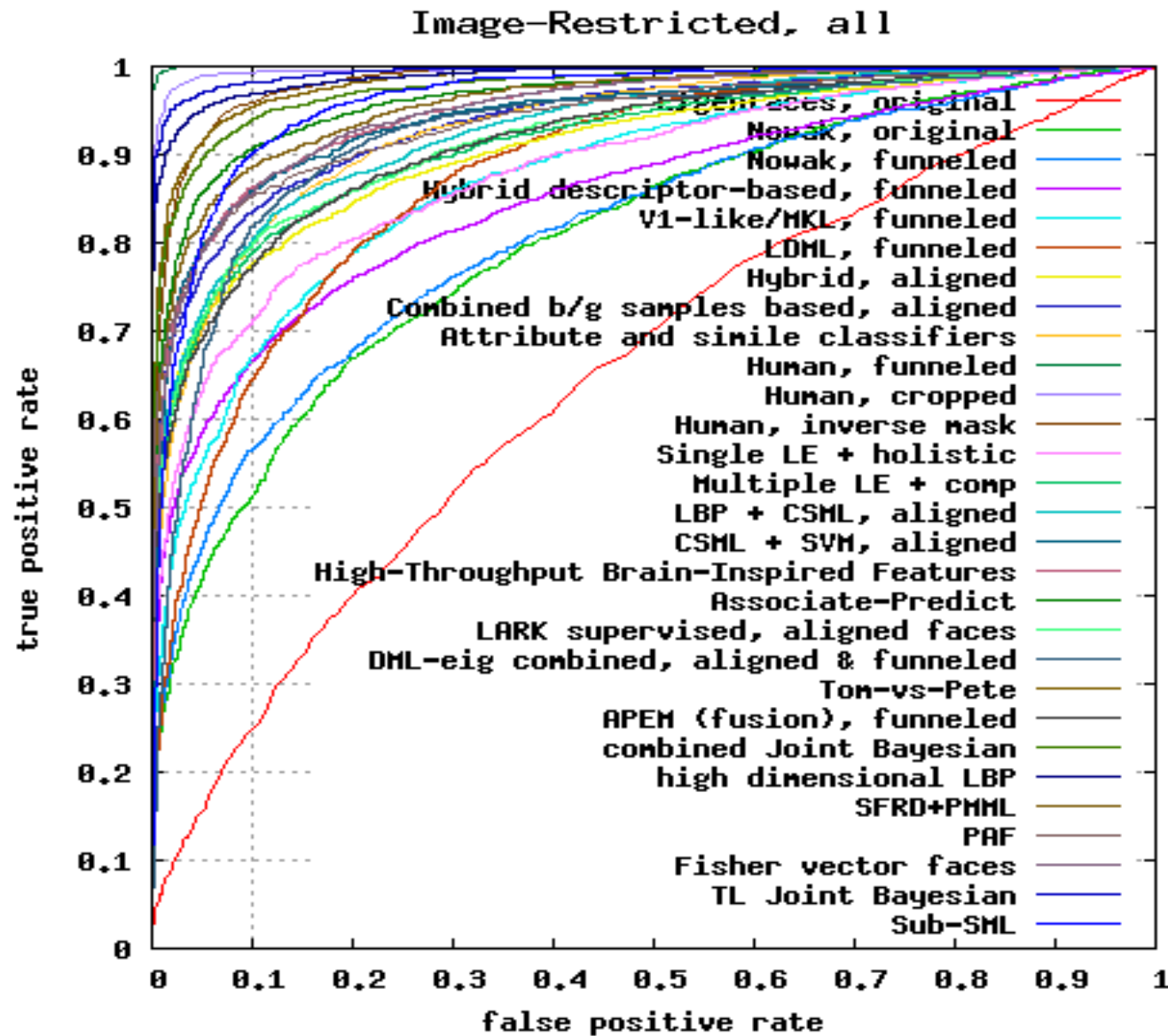
=



≠

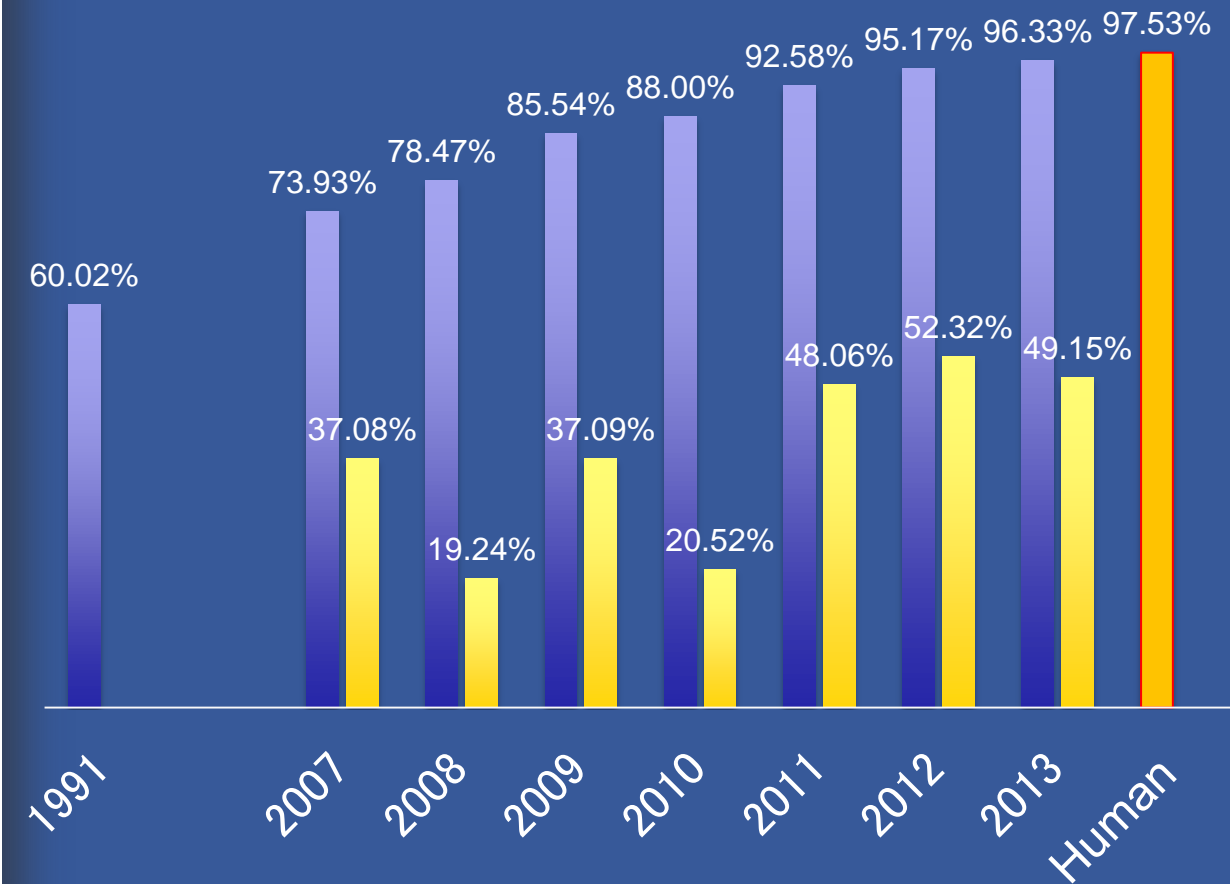


LFW: Progress over the recent 7 years

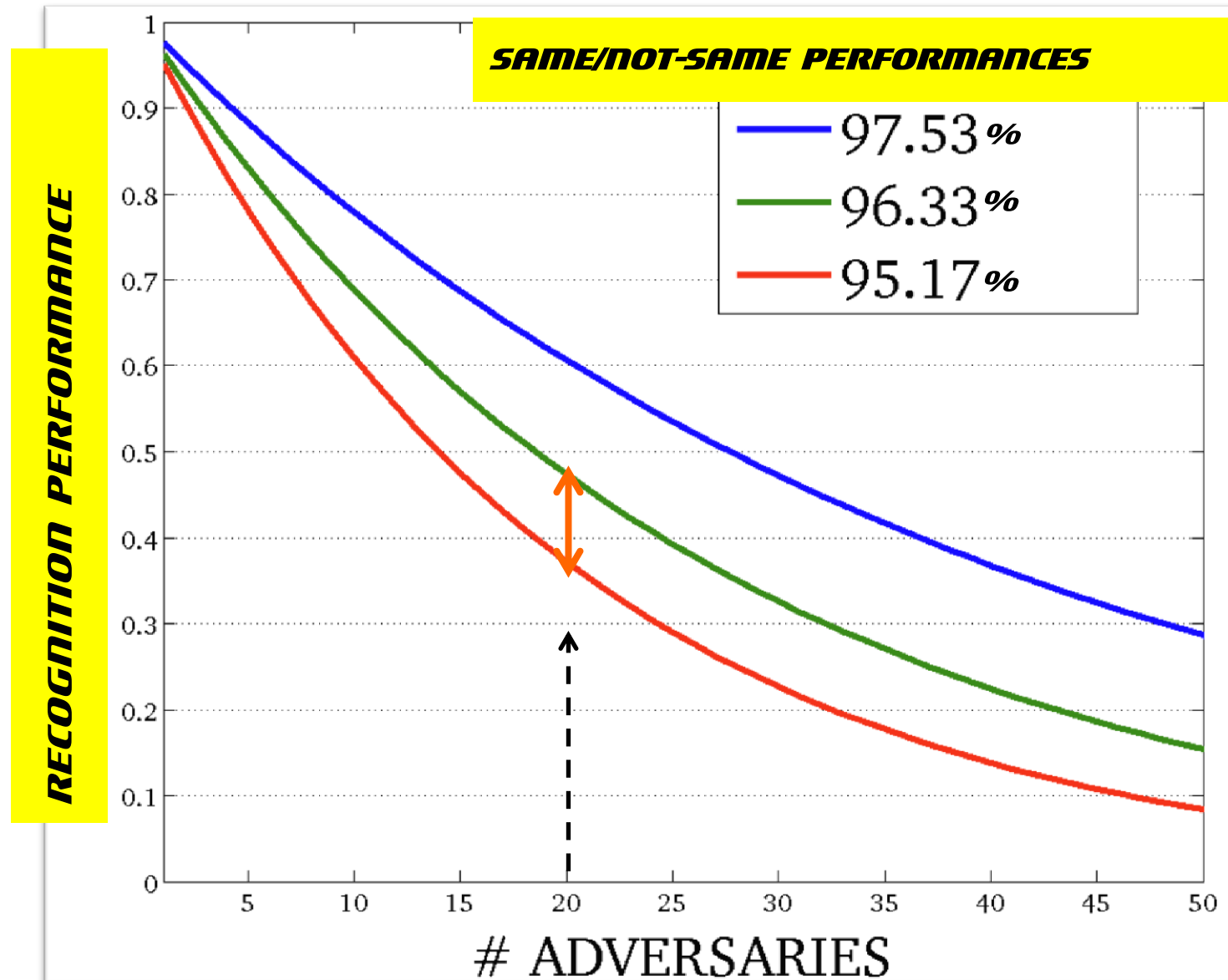


■ Accuracy / year

■ Reduction of error wrt human / year



Verification Impacts Recognition



DeepFace

*DEEPFACE: CLOSING THE GAP TO HUMAN-LEVEL PERFORMANCE IN FACE VERIFICATION;
YANIV TAIGMAN, MING YANG, MARC' AURELIO RANZATO AND LIOR WOLF (CVPR 2014)*

Face Recognition Pipeline

Detect

Align

Represent

Classify

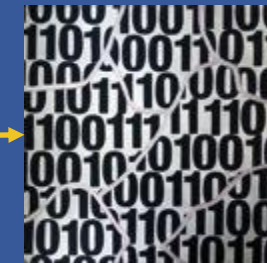
Face Recognition Pipeline

Detect

Align

Represent

Classify

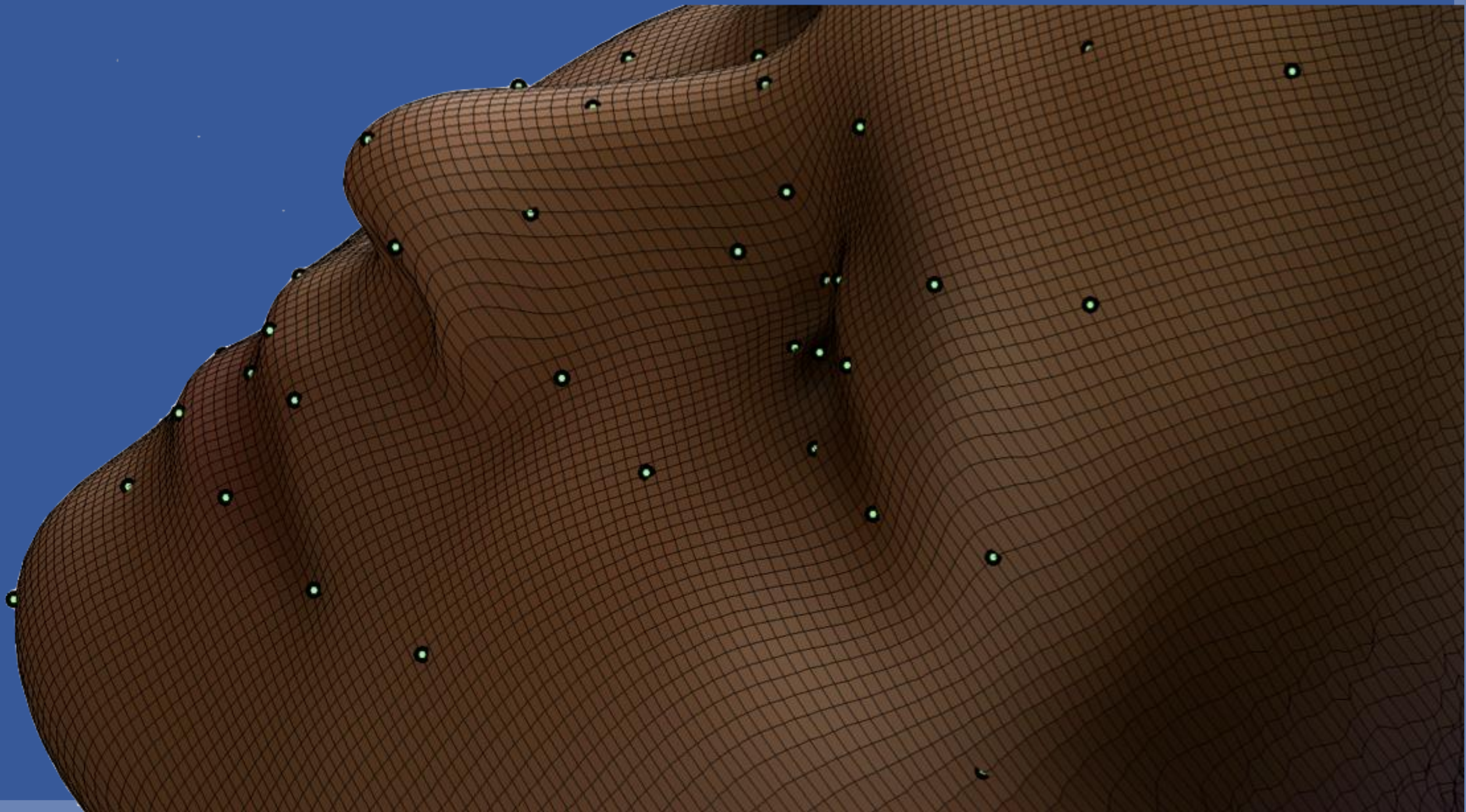


YANIV

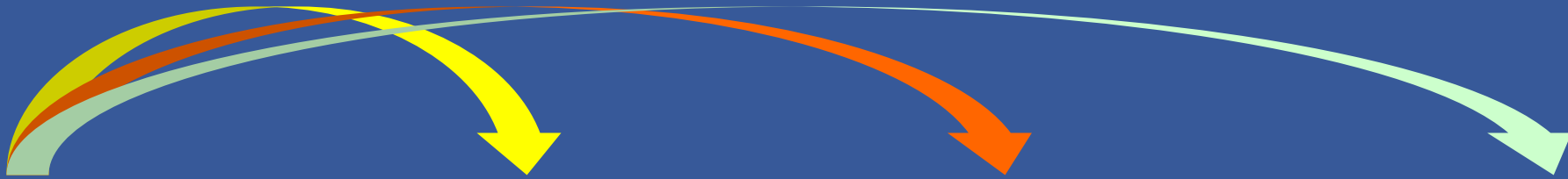
LUBOMIR

MARC AURELIO

Faces are 3D objects



Texture vs. Shape



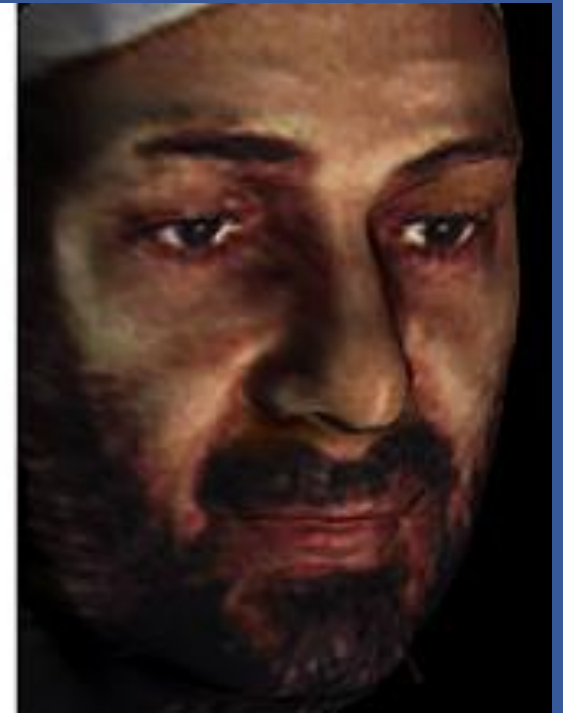
Shape A



Shape A +
Texture A



Shape A +
Texture of Bush



Shape A +
Texture of BinLaden

Face alignment

(*'Frontalization'*)



Detect

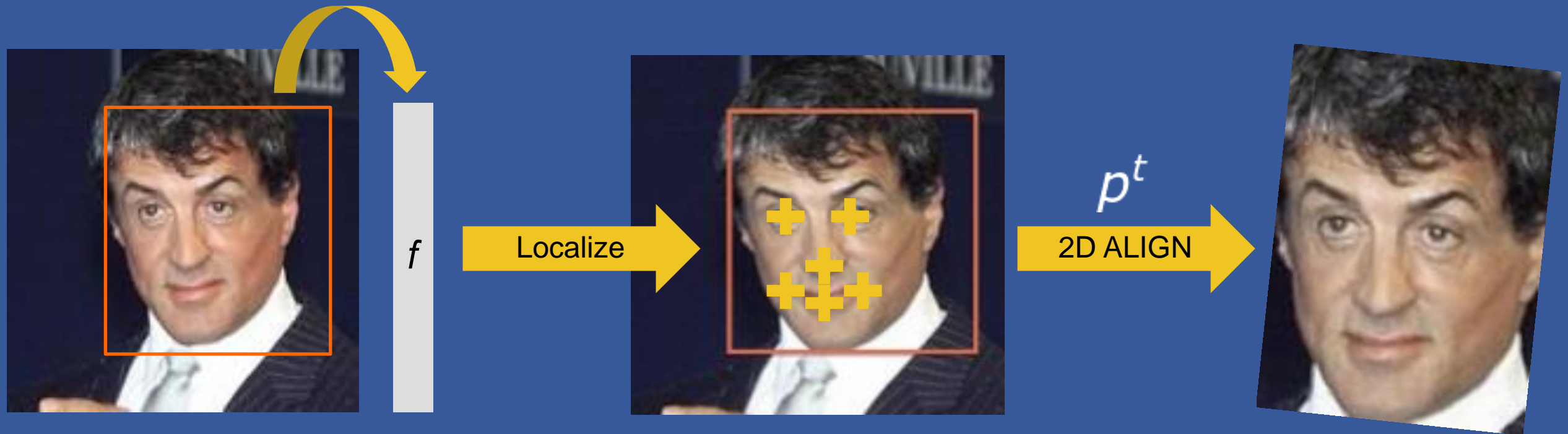


2D-Aligned



3D-Aligned

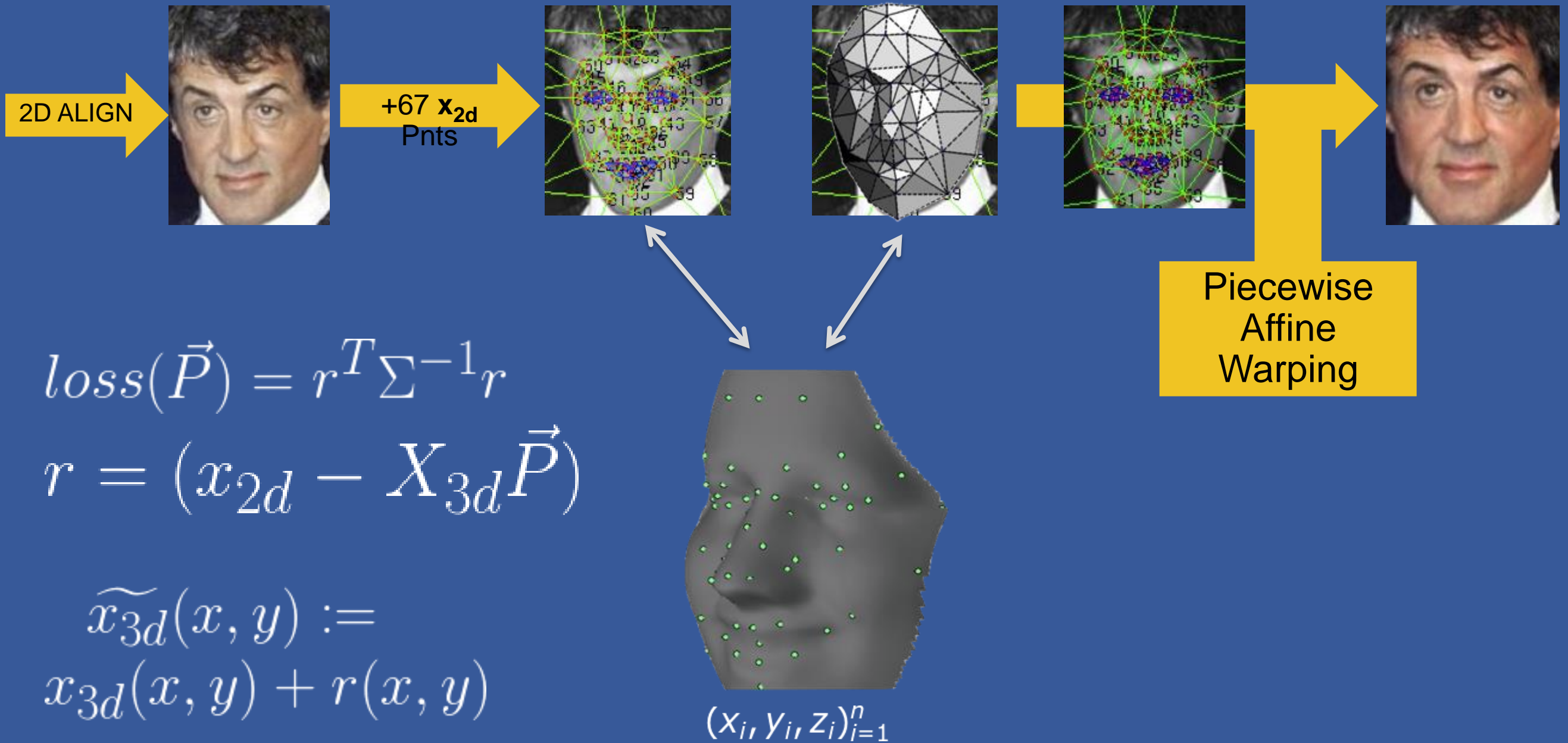
2D alignment



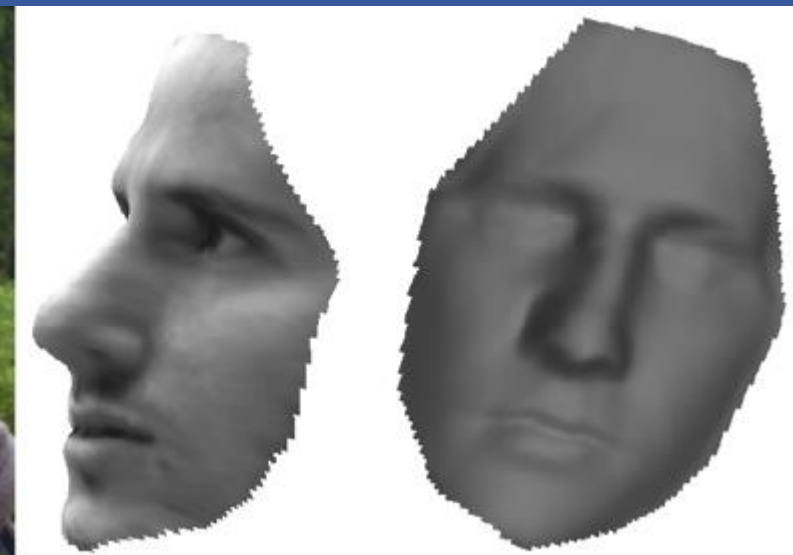
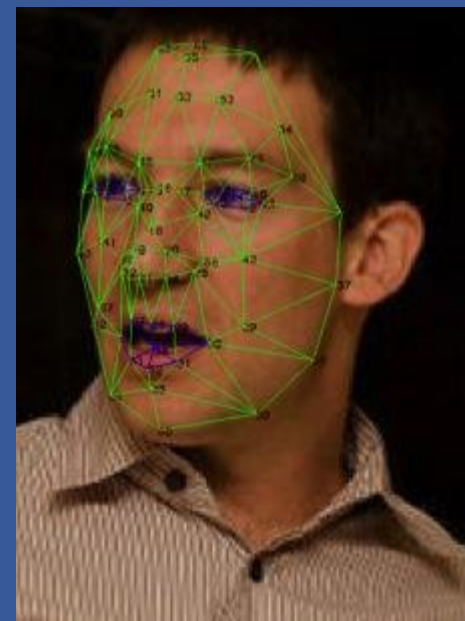
$$p^0 = A_{n \times d} \cdot f$$

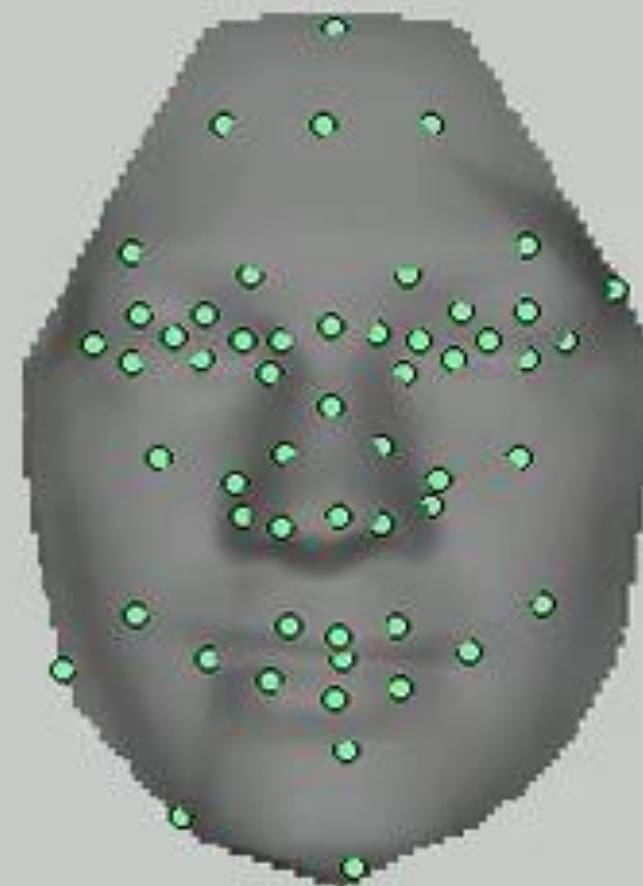
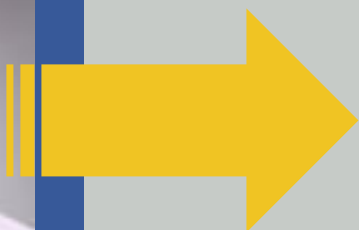
$$p^t = s_t[R_t | t_t] \cdot p^{t-1}$$

3D alignment



Examples





Next: Representation Learning

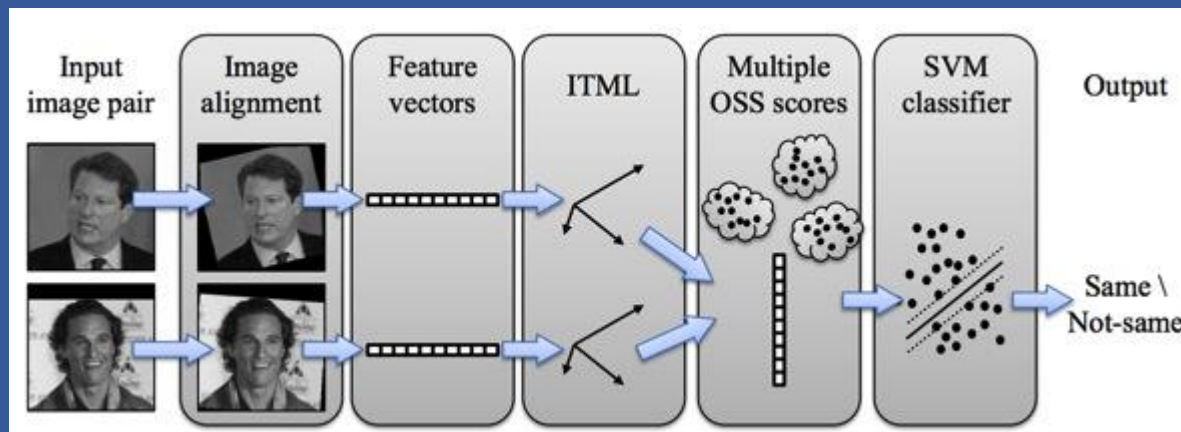
Detect

Align

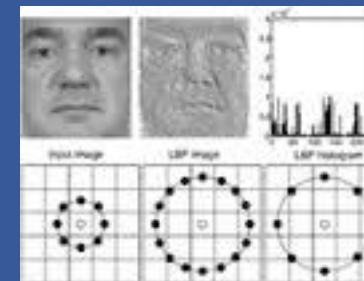
Represent

Classify

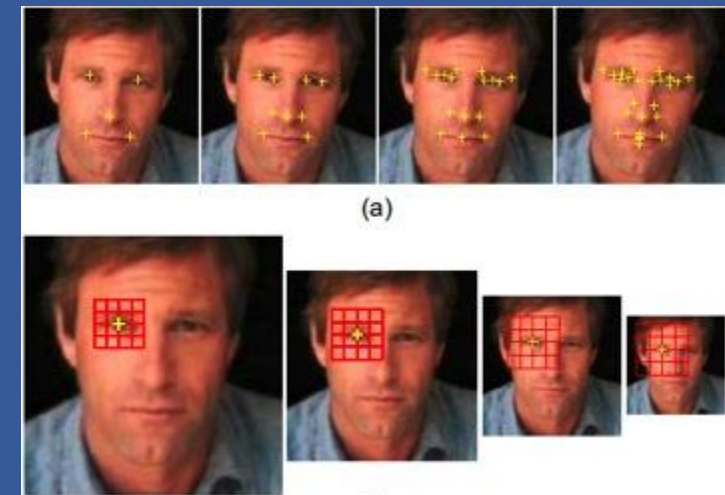
- 2004 – 2013 : Feature engineering monopoly, mostly LBP.
 - Contributions mainly in Classification.



'MULTI-SHOTS' ; TAIGMAN, HASSNER, WOLF



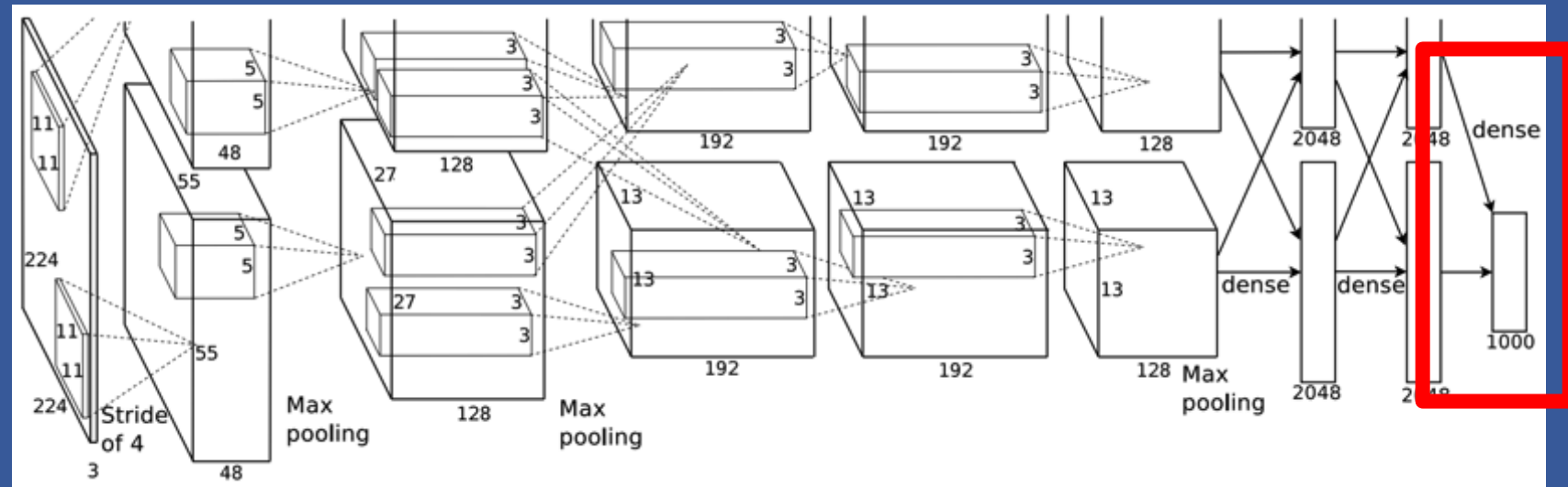
LBP; AHONEN 2004



HIGH-DIM LBP; CHEN, CAO, WEN, SUN

- 2012 : The resurrection of LeCun's Deep Convolutional Neural Networks (CNNs) by Krizhevsky, Sutskever and Hinton.

CNNs for: Image Classification vs. Face Recognition



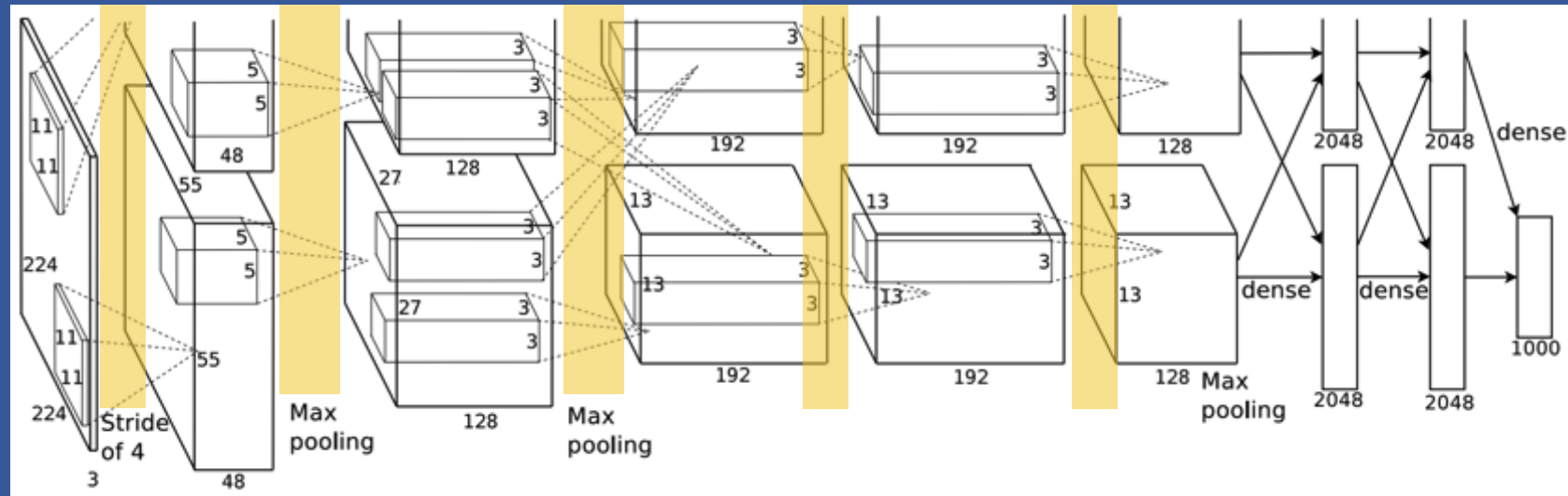
1. WE MOSTLY CARE ABOUT FEATURE LEARNING

- ***WE DO NOT KNOW THE NUMBER OF IDENTITIES BEFORE-HAND***
- ***TRANSFER LEARNING***

→ ***LAST LAYER CAN BE REMOVED OR REPLACED***

→ ***WE STILL NEED TO THINK ABOUT THE CLASSIFICATION STAGE (LATER)***

CNNs for: Image Classification vs. Face Recognition



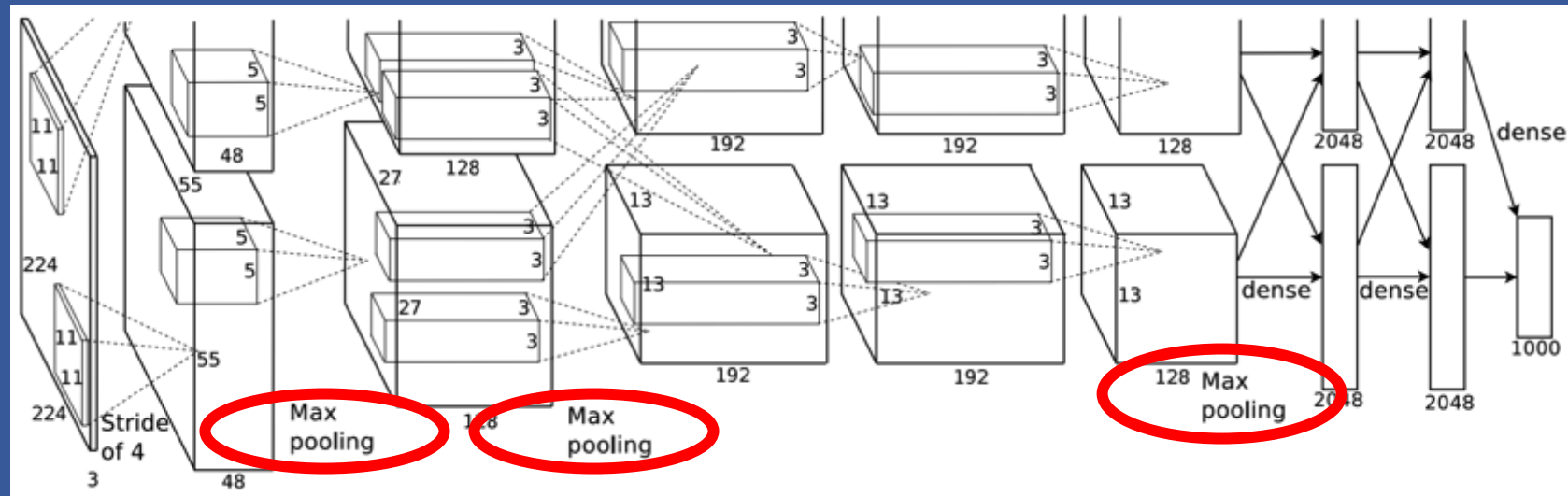
2. *GEOMETRY IS PHYSICALLY RELAXED:*

- *TRANSLATION, SCALE AND 2D-ROTATION DUE TO DETECTION AND 2D ALIGNMENT*
- *OUT-OF-PLANE ROTATION DUE TO 3D ALIGNMENT.*

*ALIGNED PIXELS → ENABLES UNTYING THE WEIGHTS → ‘**LOCALLY-CONNECTED**’ LAYERS.*

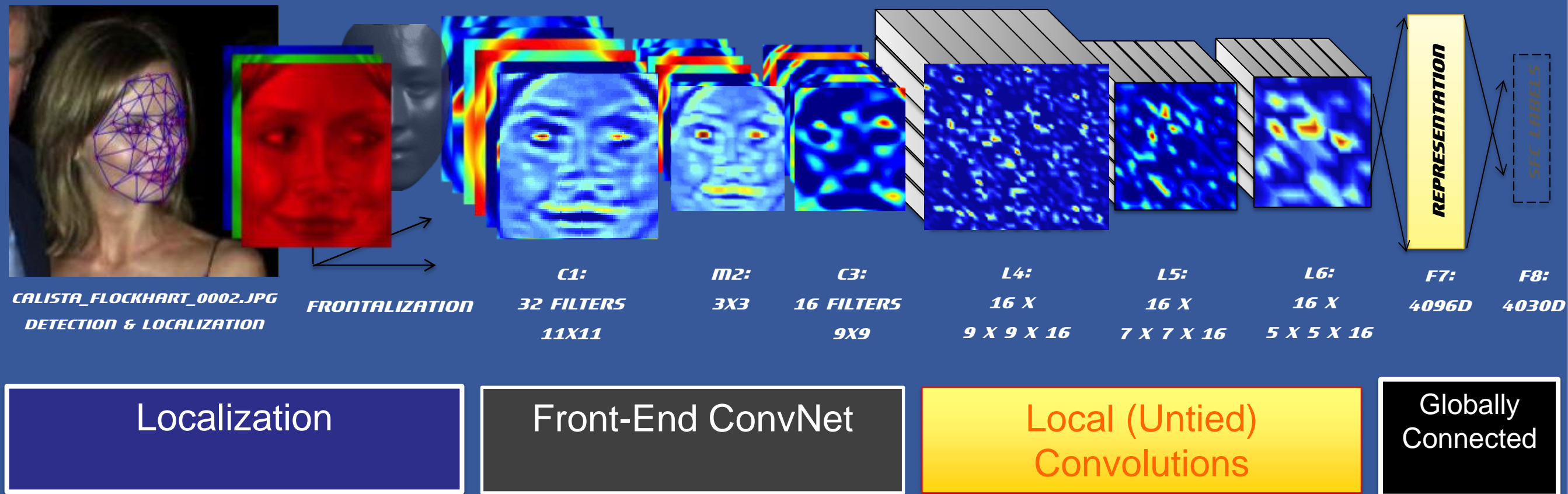
→ GREATER FOCUS IN TRAINING ON WHAT’S NOT SOLVED ALREADY.

CNNs for: Image Classification vs. Face Recognition



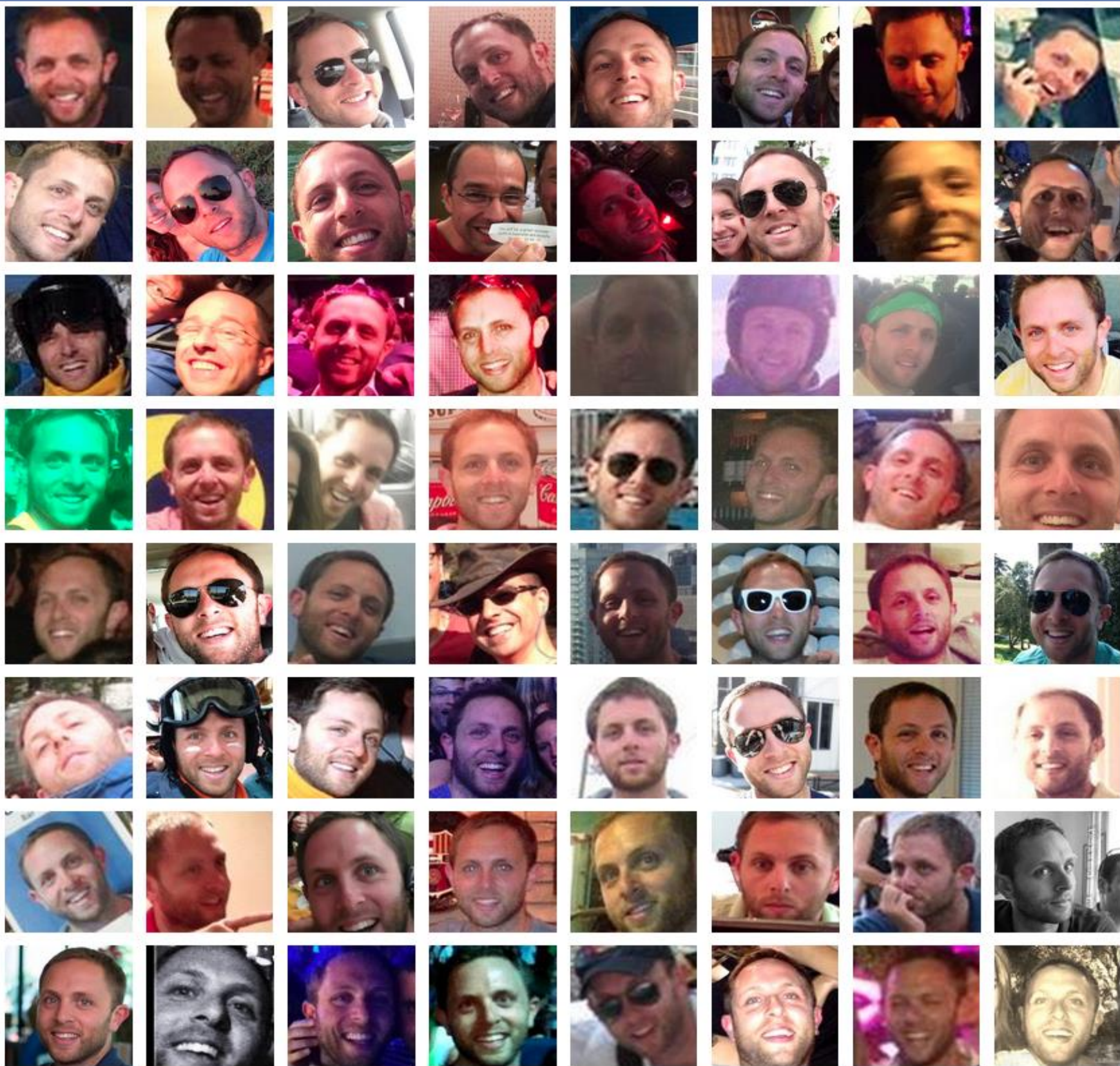
3. SEVERAL LEVELS OF (MAX-) *POOLING* WOULD CAUSE THE NETWORK TO LOSE INFORMATION ABOUT THE PRECISE POSITION OF DETAILED FACIAL STRUCTURE AND MICRO-TEXTURES.

DeepFace Architecture



$$G(I) = g_{\phi}^{F_7}(g_{\phi}^{L_6}(\dots g_{\phi}^{C_1}(T(I, \theta_T)) \dots))$$

ALIGNMENT



SFC TRAINING

DATASET

(PRE-CROPPING)

***4.4 MILLION PHOTOS
BLINDLY SAMPLED,
CONTAINING MORE
THAN 4,000
IDENTITIES
(PERMISSION
GRANTED)***

Detect

Align

Represent

Classify



DeepFace
Replica



DeepFace
Replica

(A) COSINE ANGLE

$$S(f_1, f_2) = \frac{\langle f_1, f_2 \rangle}{\|f_1\| \|f_2\|}$$

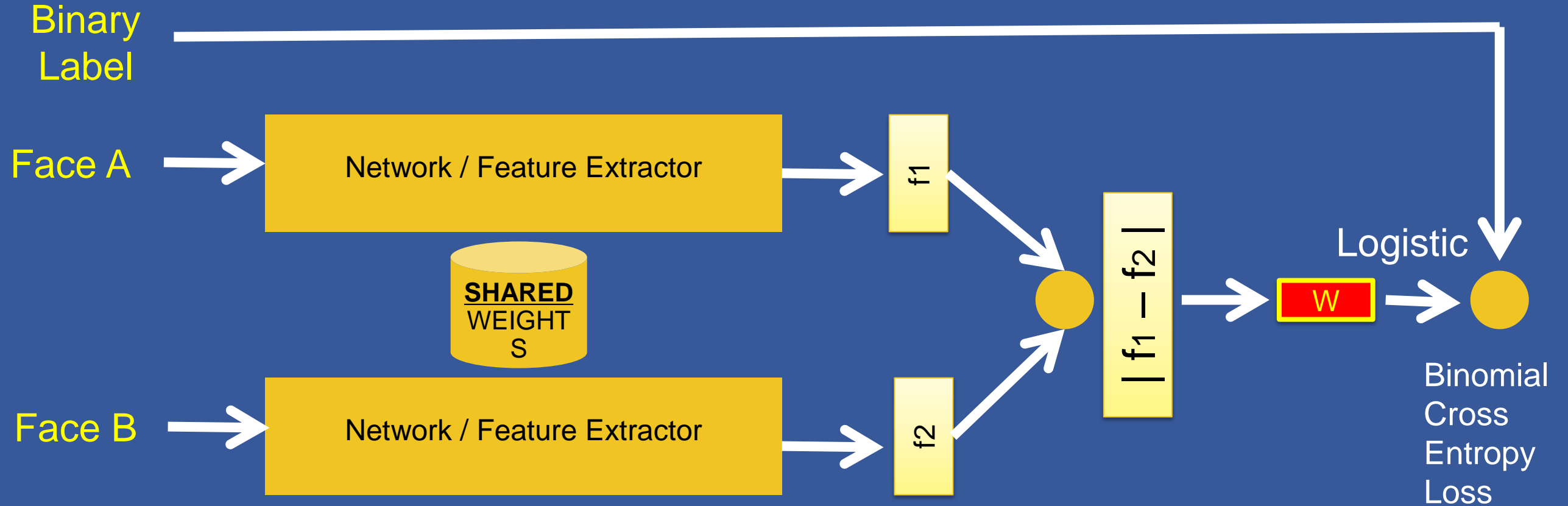
(B) KERNEL METHODS

$$S_{\chi^2}(f_1, f_2) = \sum w_i \frac{(f_1[i] - f_2[i])^2}{f_1[i] + f_2[i]}$$

(C) SIAMESE NETWORK¹

$$S_{Siam}(I_1, I_2) = \frac{1}{1 + e^{-(W|f(I_1) - f(I_2)| + b)}}$$

Deep Siamese Architecture [1]

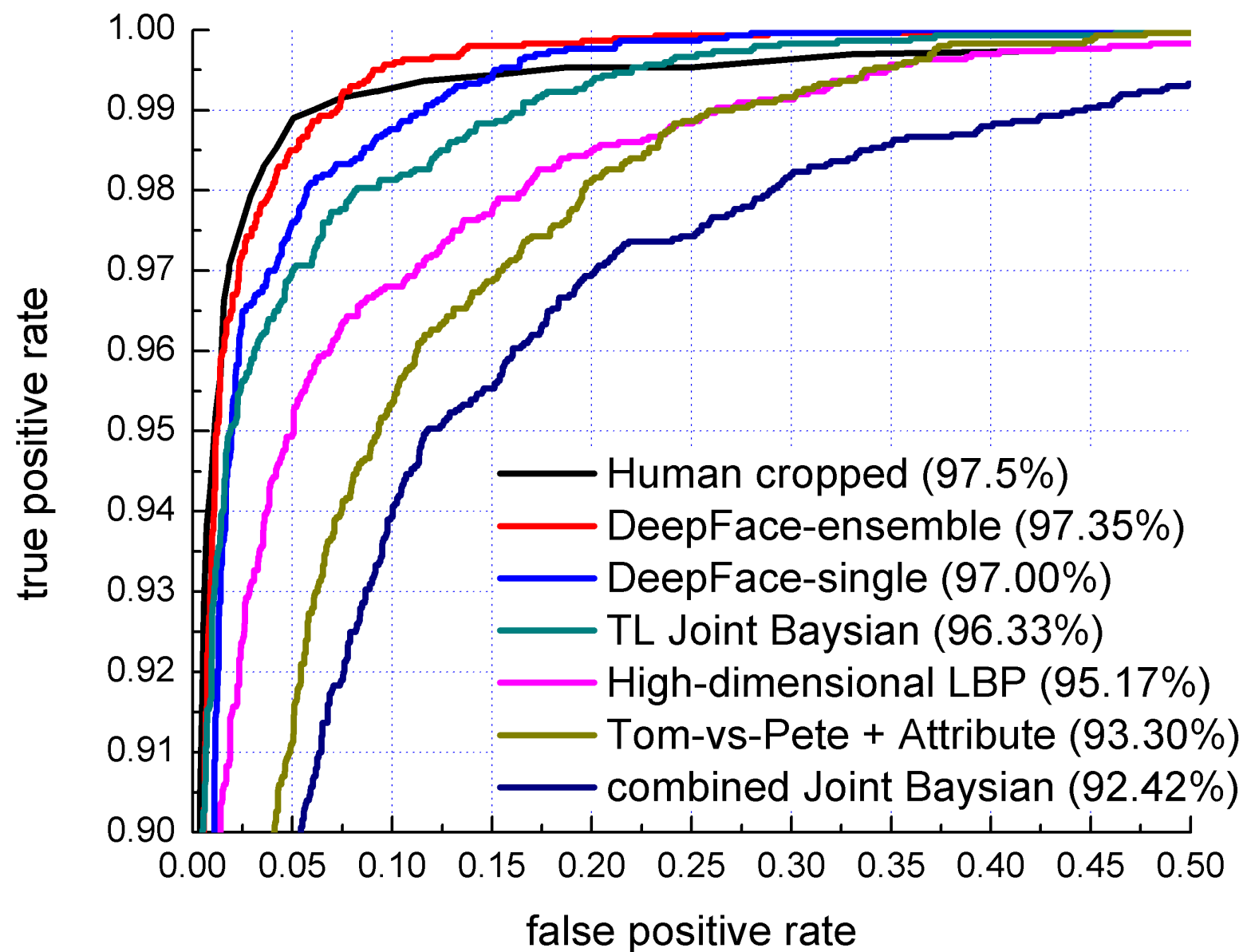


$$p = \frac{1}{1 + e^{-(W * |f_1 - f_2| + b)}}$$

$$E = -y \log(p) - (1 - y) \log(1 - p)$$

[1] Dimensionality Reduction by Learning an Invariant Mapping - Hadsell, Chopra, LeCun (2006)

Results on LFW



‘Explaining’ the False Negatives pairs (1.65%)



age

sunglasses

occlusion/
hats

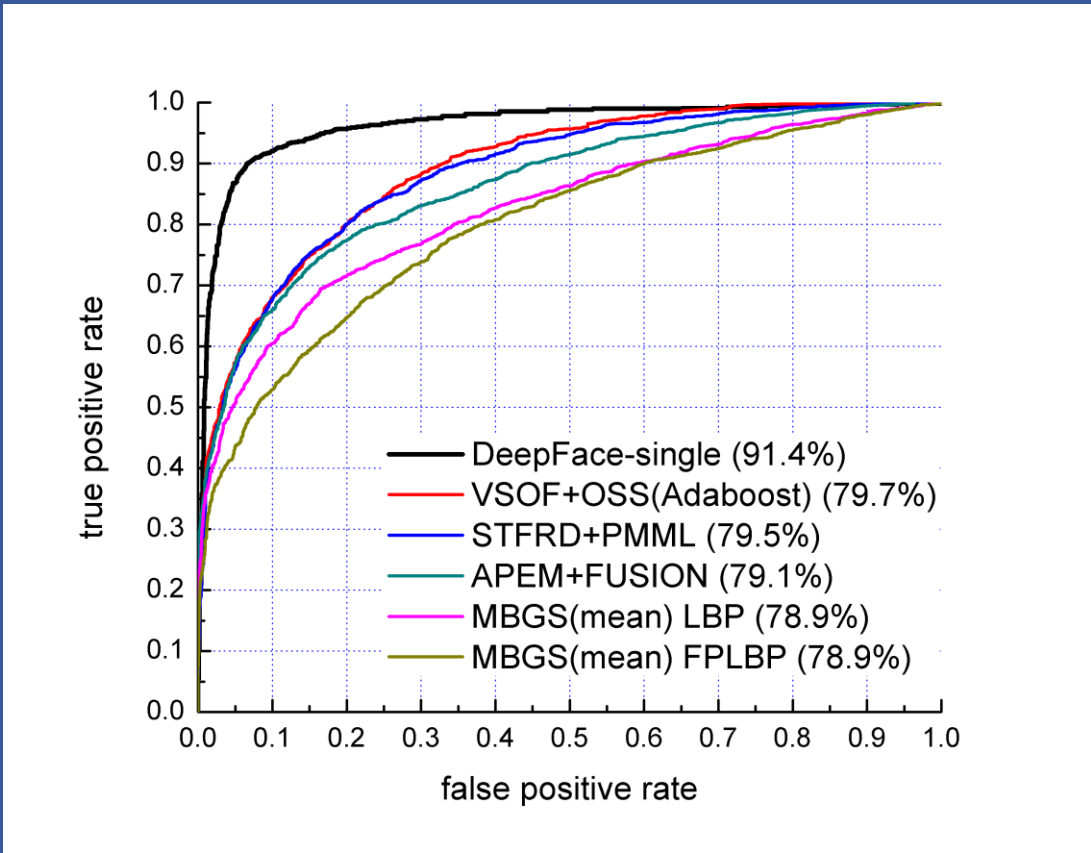
profile

errata

False Positive pairs (1.00%)



Results on YouTube Faces (Video)



↑
False negatives

False positives
→



Face Identification (1:N)



Unaccounted challenges in **verification**:

- I. Reliability
- II. Large confusion ($P \times G$)
- III. Different distributions
- IV. Unknown class



LFW Identification (1:N) Protocols²

1. Close Set

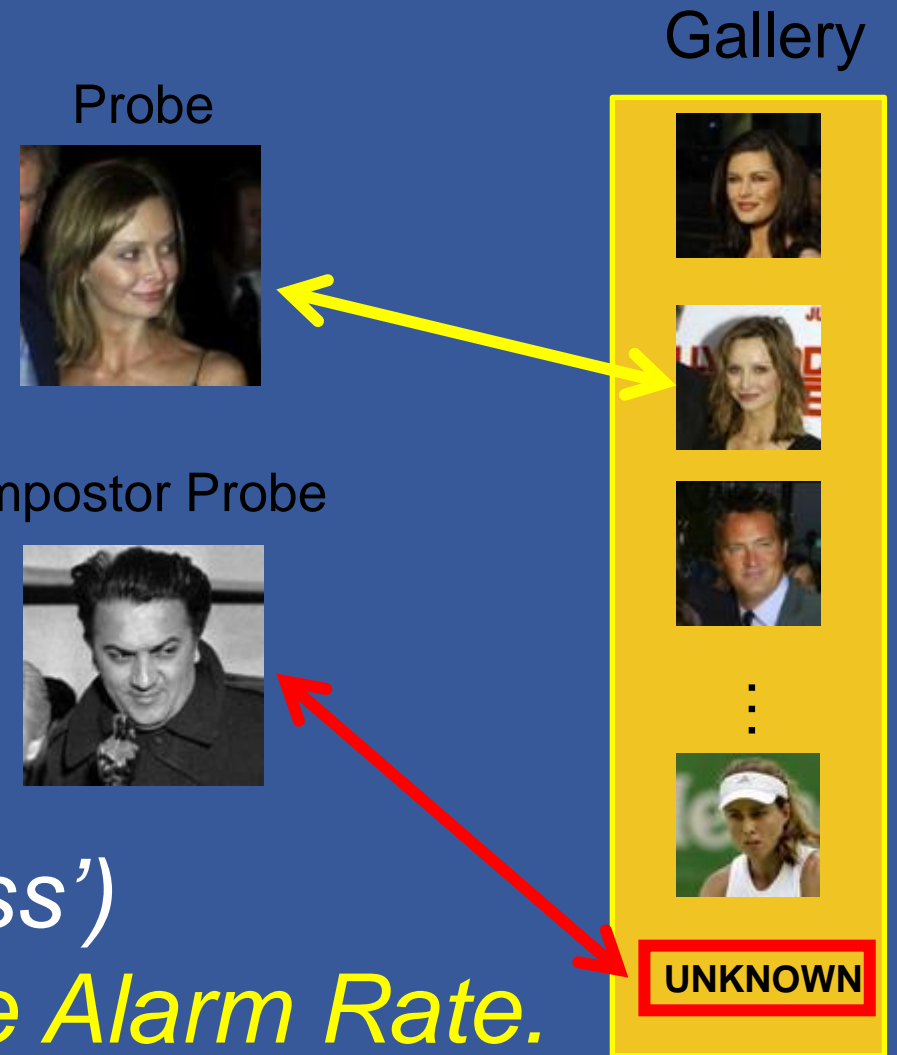
- #Gallery¹: 4,249
- #Probes: 3,143

Measured³ by *Rank-1 rate*.

2. Open Set

- #Gallery¹: 596
- #Probes: 596
- #Impostors: 9,491 ('unknown class')

Measured³ by *Rank-1 rate @ 1% False Alarm Rate*.



¹ Each identity with a **single** example

² Unconstrained Face Recognition: Identifying a Person of Interest from a Media Collection
Best-Rowden, Han, Otto, Klare and Jain (Technical Report MSU-CSE-2014-1)

³ Training is **not** permitted on LFW ('unsupervised')

LFW Identification (1:N) Results

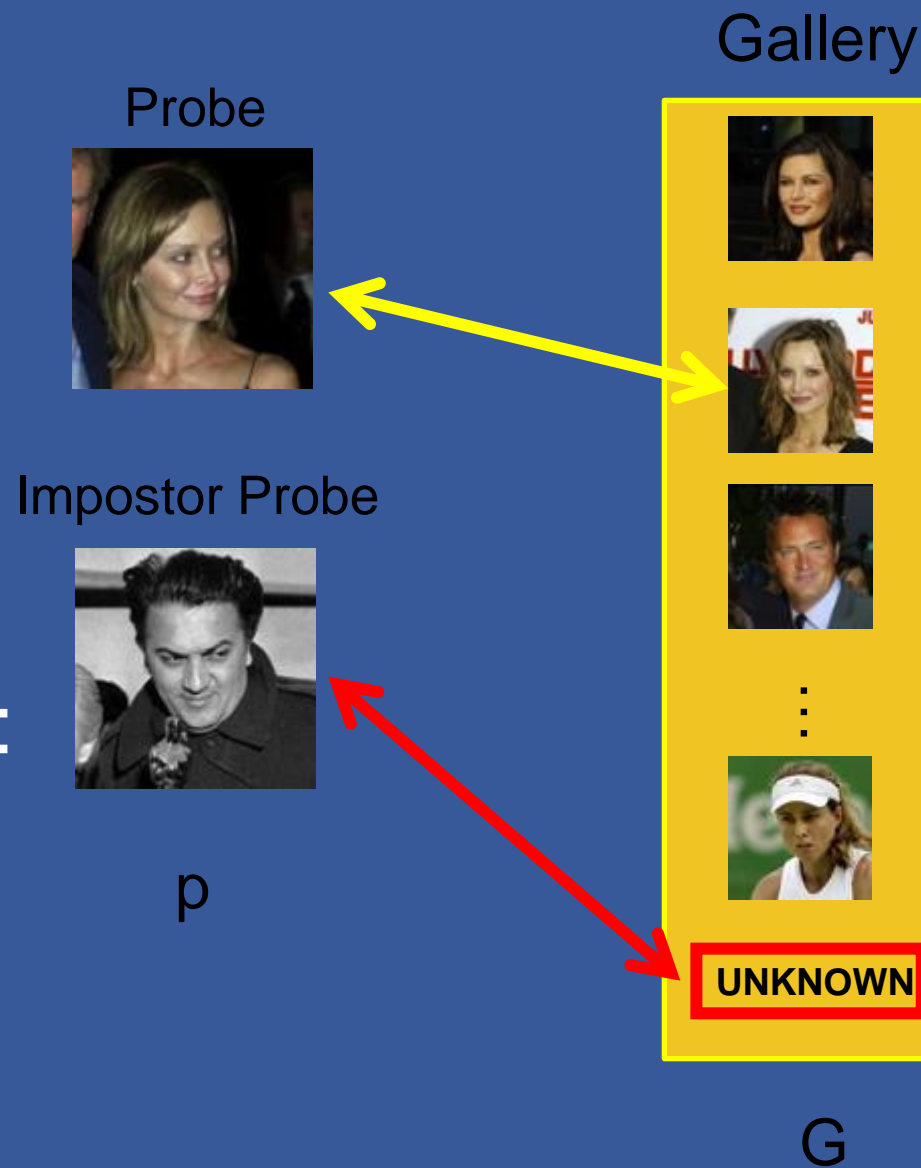
Method	DeepFace [20]	BLS [3]*	NIST's s1 [1]
Verification	97.35	93.18	-
Rank-1	64.9	18.1	56.7
DIR @ 1%	44.5	7.89	25

Cosine similarity measure ('unsupervised') :

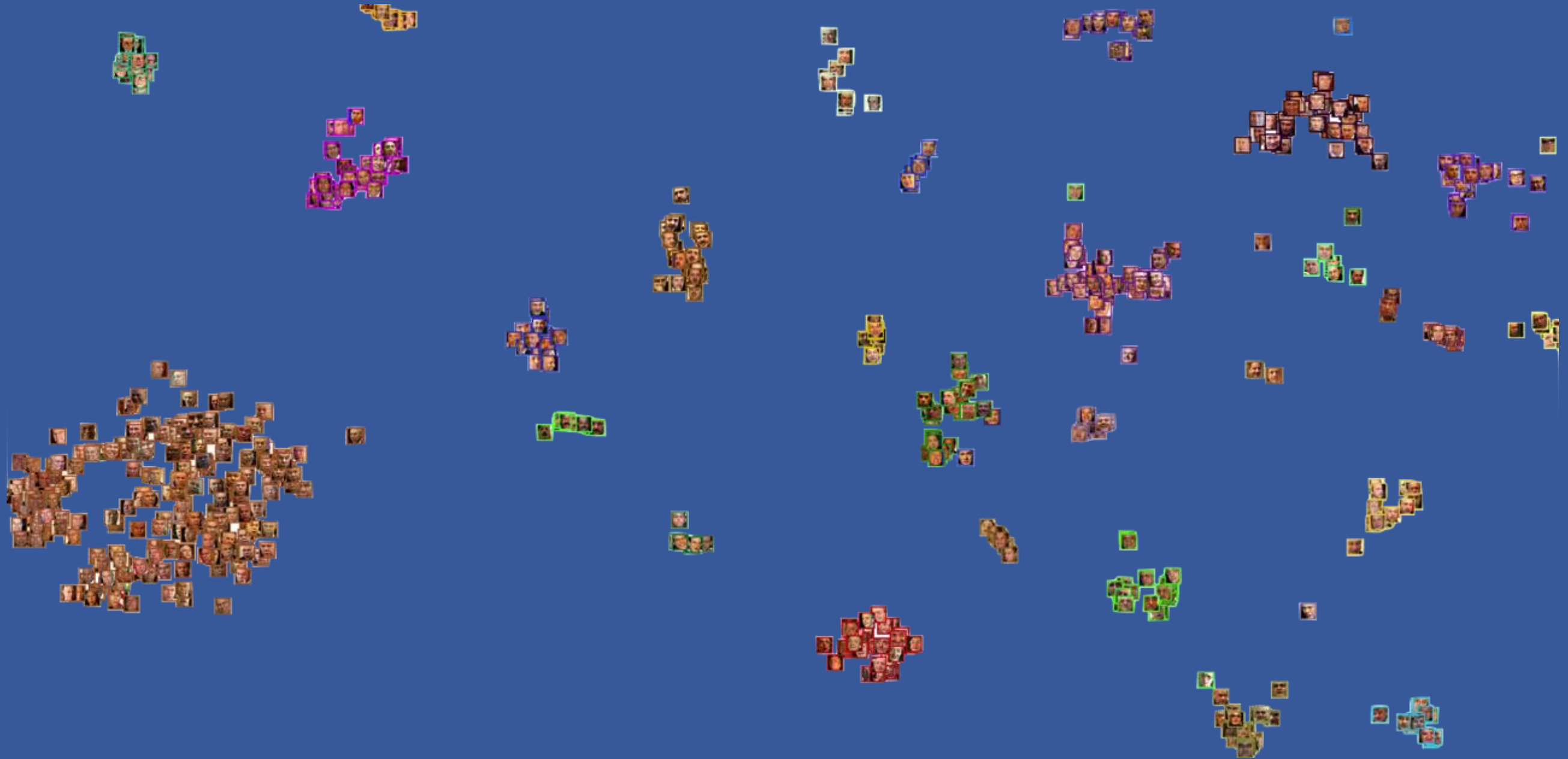
$$\text{Confusion Matrix} = G^T * P$$

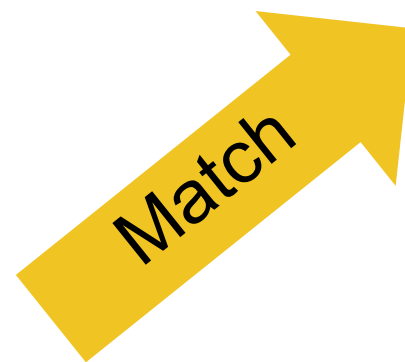
G is 256 x 4249

P is 256 x 3143



(part of the) t-SNE visualization of LFW faces





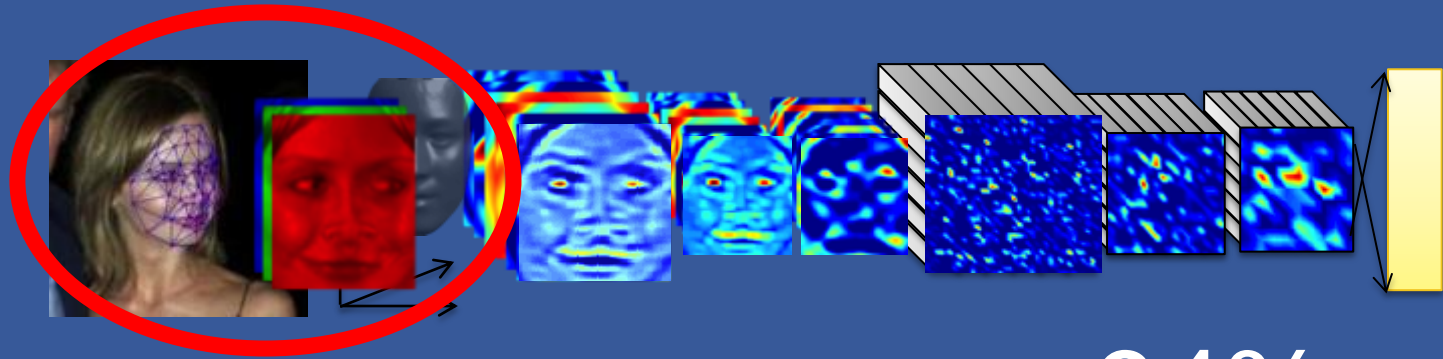
Why does it work so well ?

1. **Coupling** alignment with Locally-Connected layers
2. Large capacity model that actually enjoy large data

But can we understand more with respect to the roles of:

- What each layer is actually doing
- Is alignment necessary?
- Is regularization needed?
- Dimensionality & Sparsity
- Will more data help?

Localization is needed



75%

Original
+ ImageNet

89%

Original

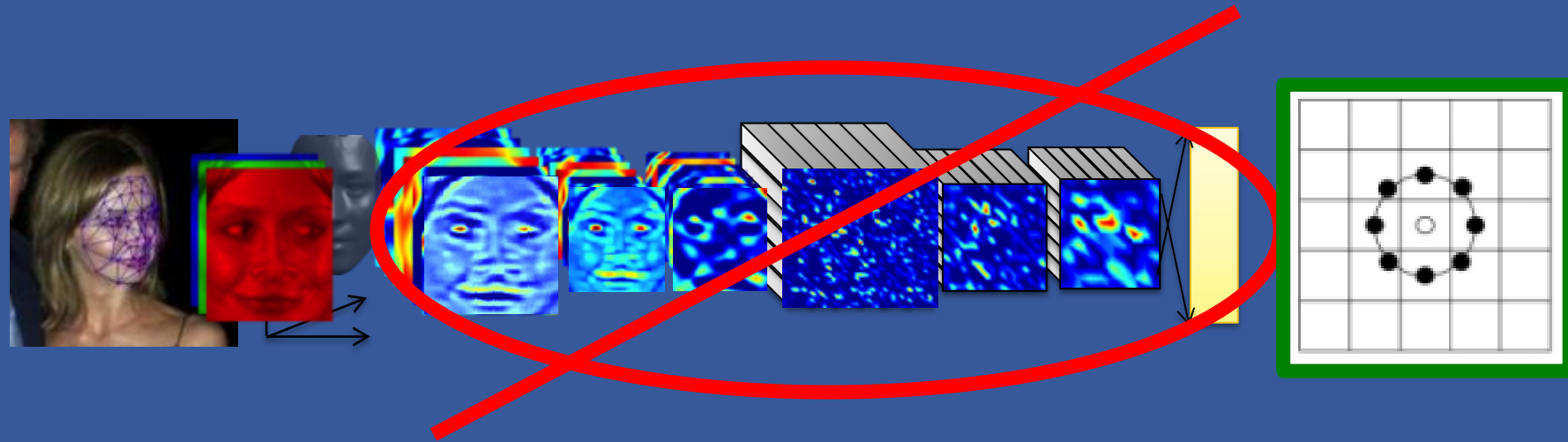
94%

2D-Aligned

97%

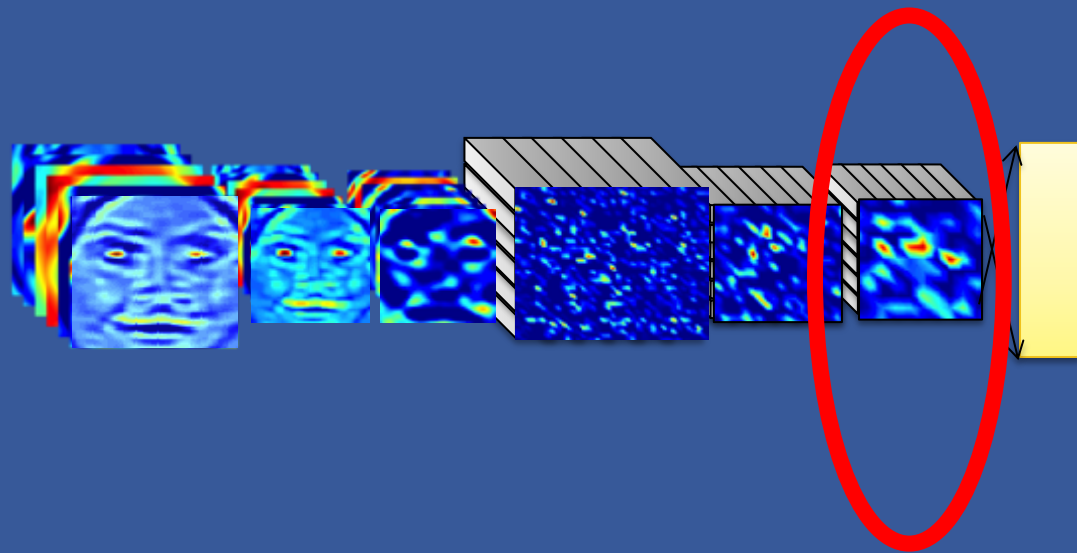
3D-Aligned

Localization is needed but insufficient

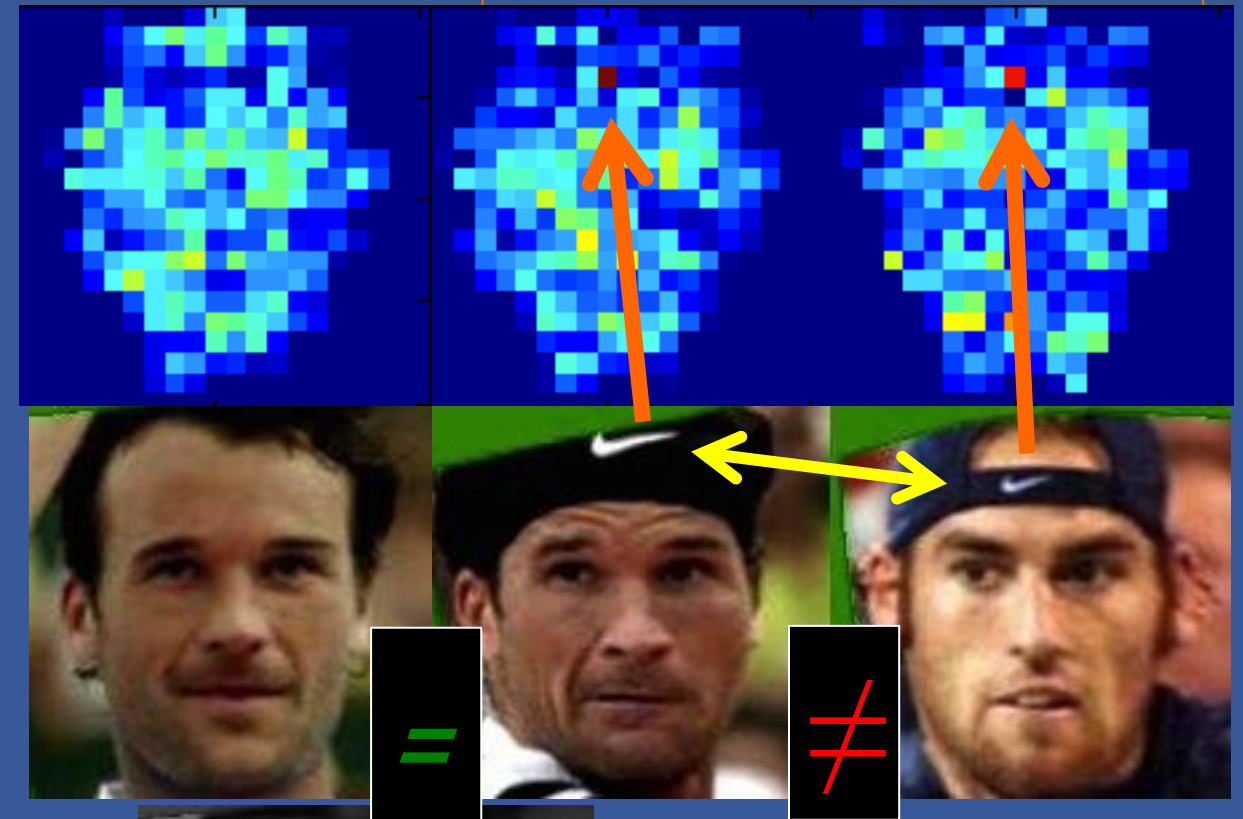


- Alignment – DNN + LBP → Accuracy drops to 91.5% (-6%)

Local Patches are Insufficient

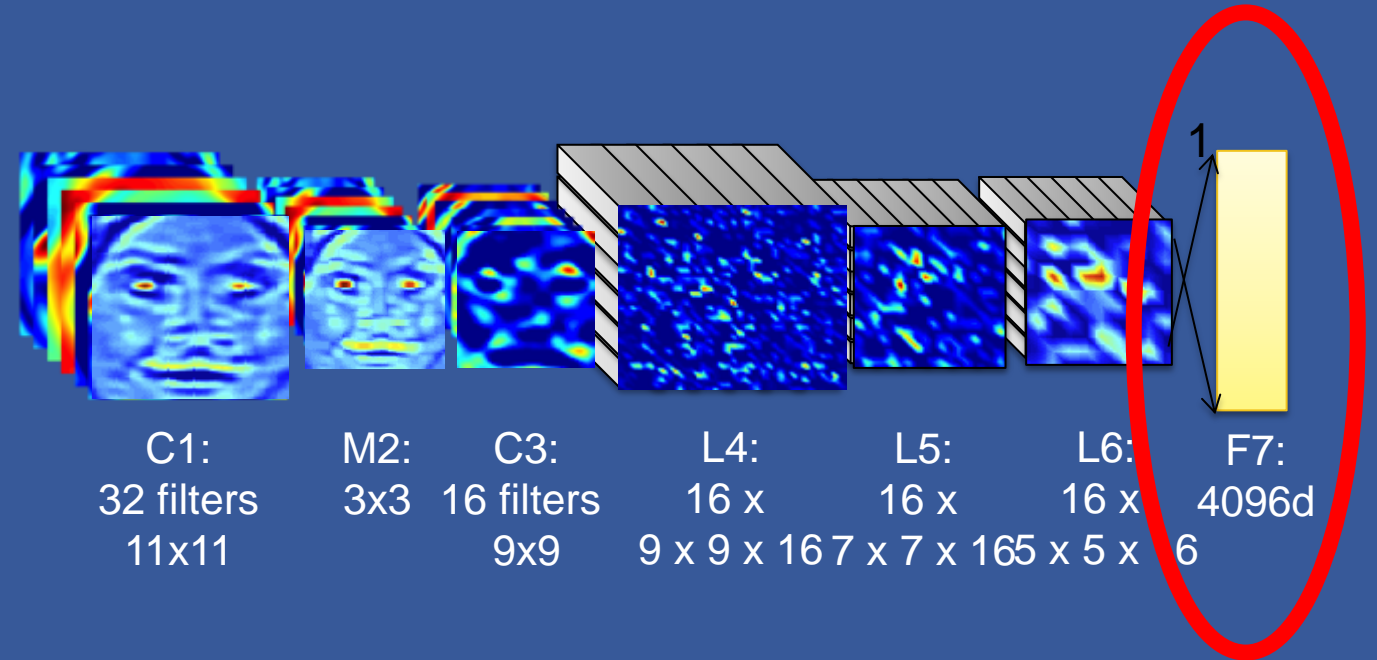


FALSE POSITIVE



→ Fully-Connected Layer is the holistic representation

Projects input 'features'
Into the representation.



1. **Correlates** between different local parts
2. Can exploit **symmetries** in faces
3. **High-Level templates**, a-la Eigenfaces (PCA)

Sparsity

- The $\text{RELU} := \max(0, x)$ encourage sparsity.
- Weights can be 'thought of' as weak template classifiers:

$$\text{Output} := \max (0, W * \text{input} + \mathbf{b})$$

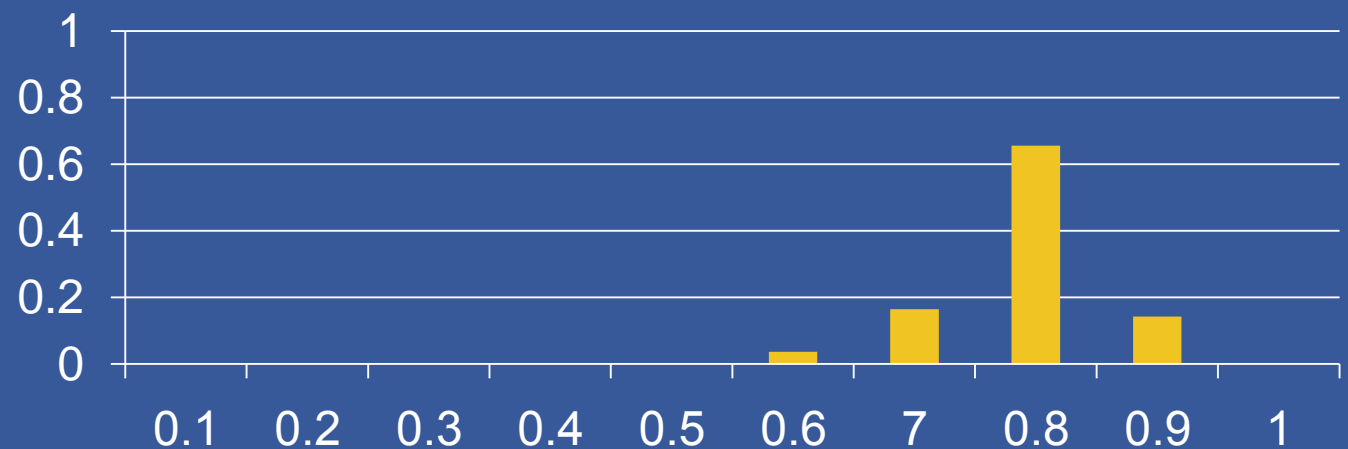
- Bias 'b' is a trainable **thresholder** / filter:

IF : $W * \text{input} < -\mathbf{b}$ THEN

 Output := 0

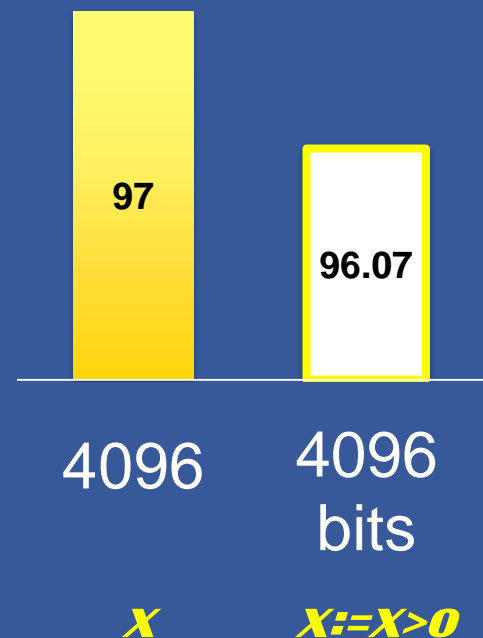
ELSE

 Output := $W * \text{input} + b$



Most of the information is encoded in whether a unit is fired or not

$X := (X > 0) \rightarrow$ Performance drops only a bit.



The **norm** of the representation is a measure of signal acquisition

FOR FACES: $\|F(i)\|$ IS A MEASURE OF FEED-FORWARD CONFIDENCE

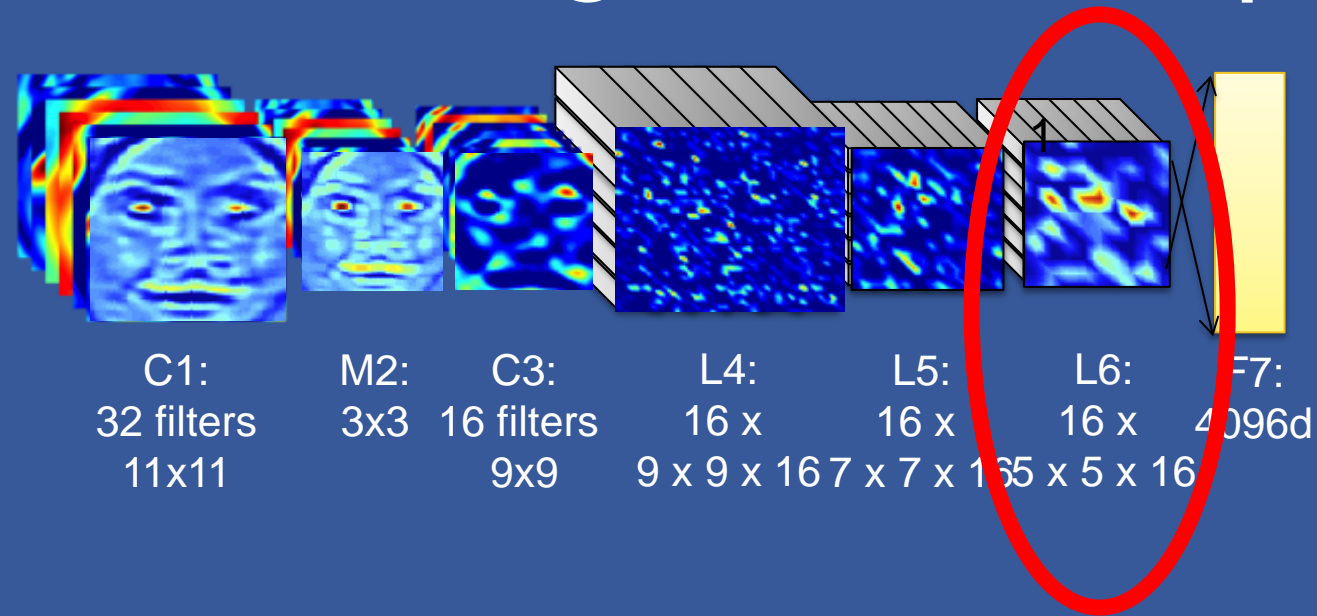
SMALLEST NORM'S IN LFW:



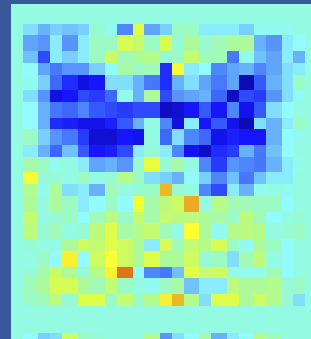
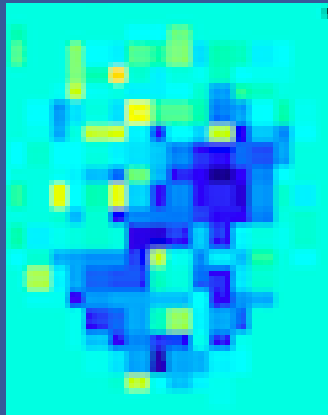
LARGEST NORM'S IN LFW:



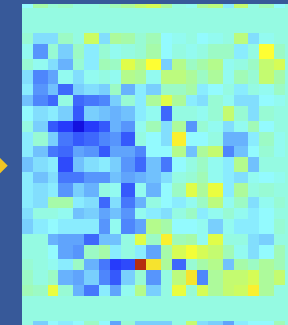
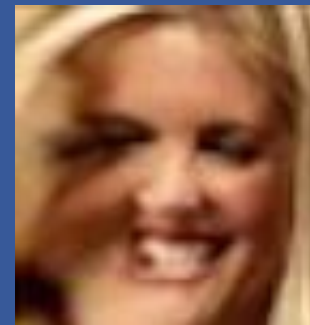
Understanding feature response



Occlusion



Failed Alignment



Correlation between norm & accuracy confidence

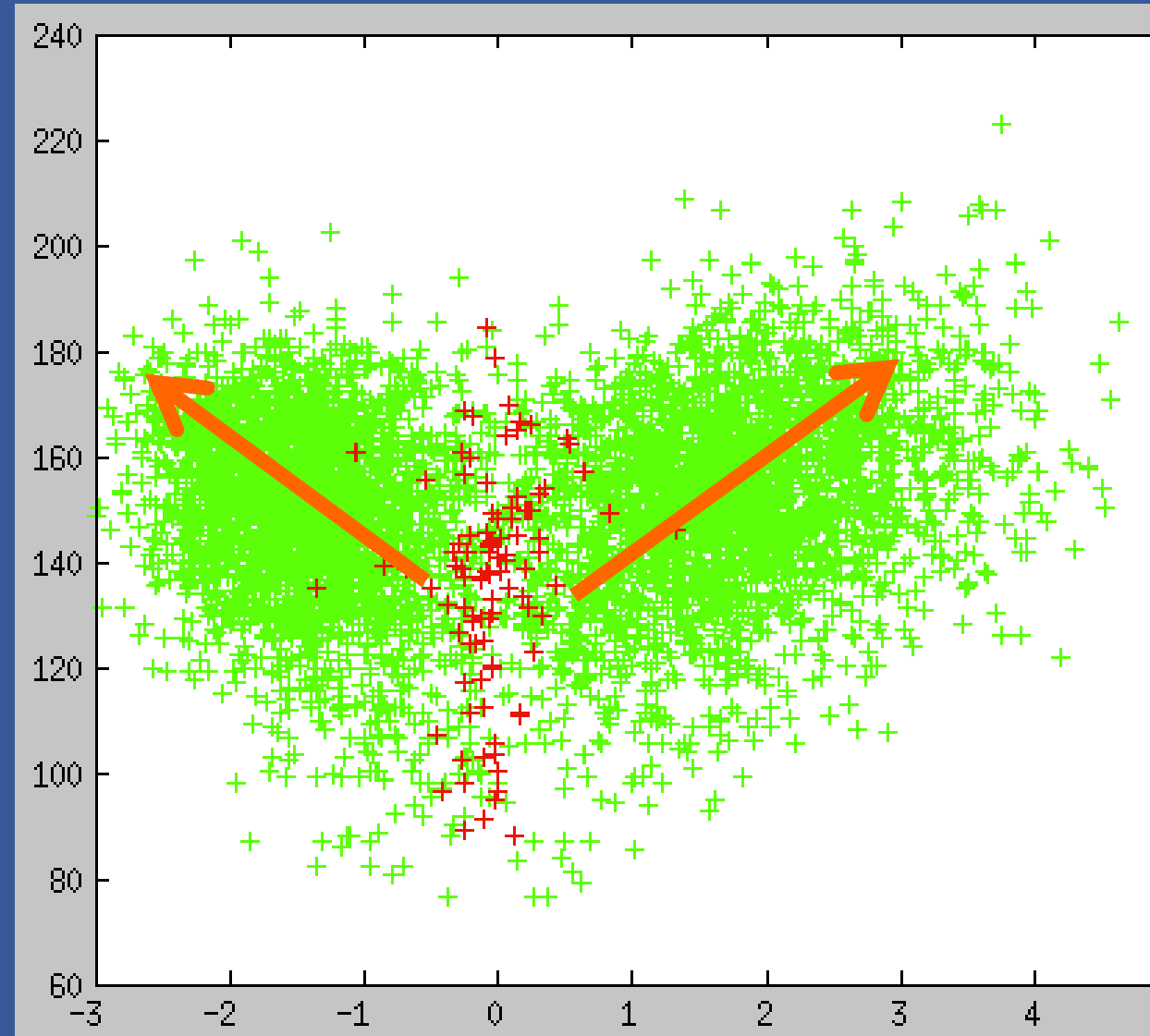


TRUE POSITIVE OR TRUE NEGATIVES



FALSE POSITIVES OR FALSE NEGATIVES

REPRESENTATION NORM

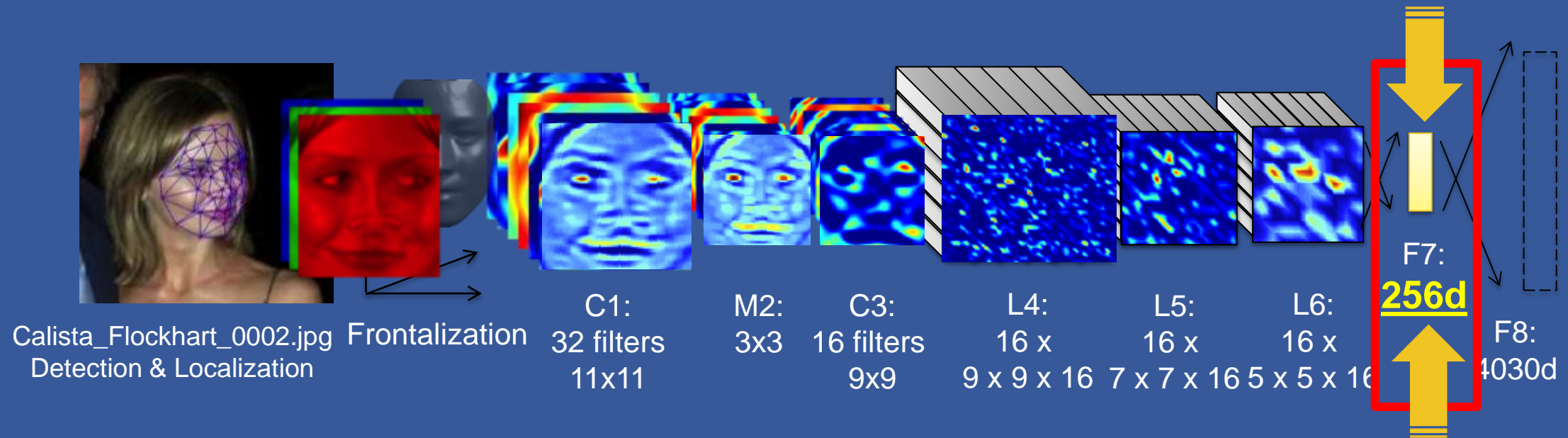


**SURELY
NO**

UNCERTAIN

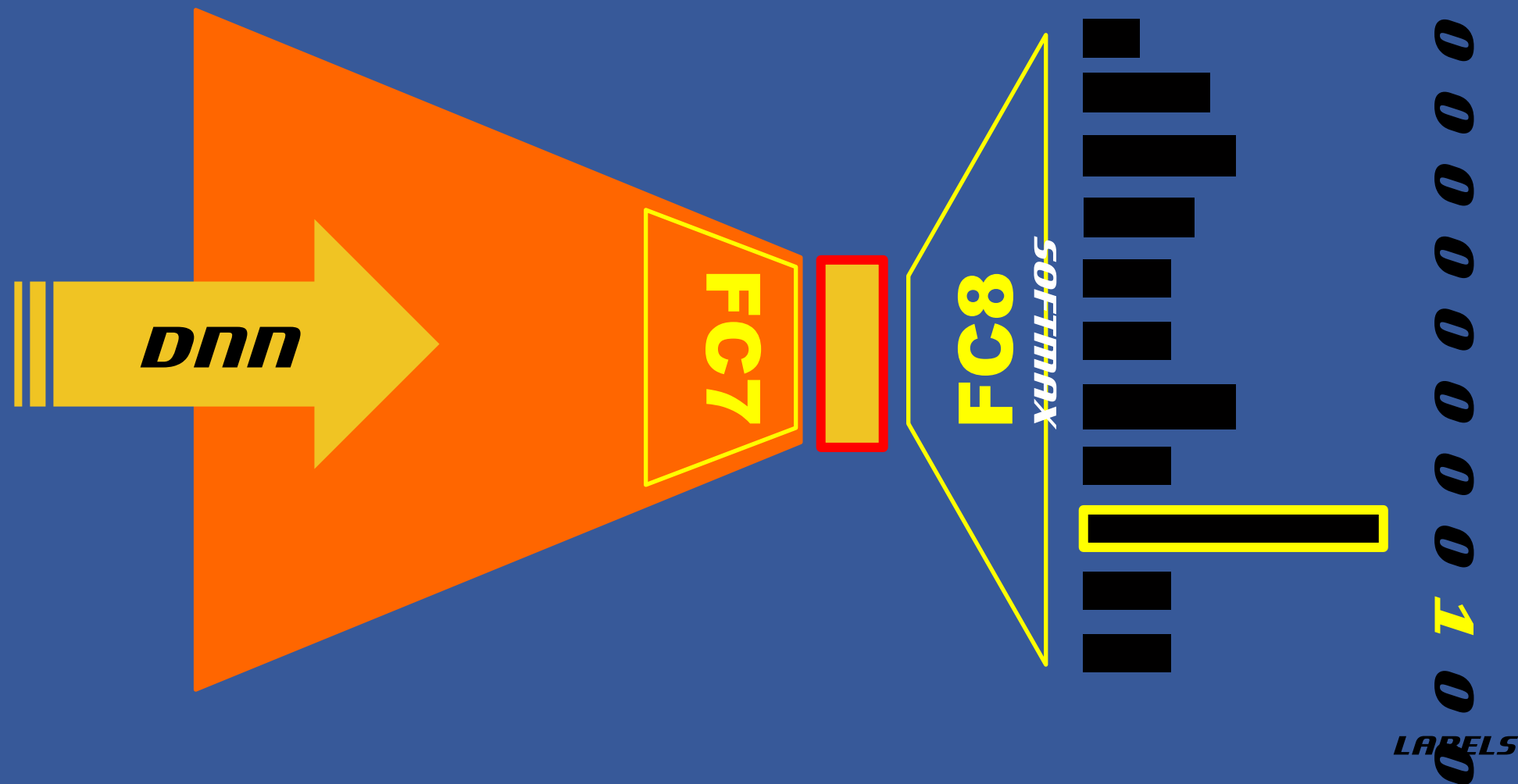
**SURELY
YES**

Bottleneck is an important Regularizer in Transfer Learning



The network overfits less on the SOURCE training set, and performs better on the TARGET when reducing the representation layer (F7) from 4K dims to **256 dims**.

Bottleneck regularizes Transfer Learning



CNN's (can) saturate

“Results can be improved simply by waiting for faster GPUs and bigger datasets to become available” -- Krizhevsky et al.

What happens when the network is fixed & the number of training grows from 4m \rightarrow 0.5b ?

Answer: our findings reveal that this holds to a certain degree only.

Data is practically infinite.



The Flickr logo, consisting of the word "flickr" in a bold, lowercase sans-serif font. The letters "f", "l", "i", "c", "k", and "r" are blue, while the letter "r" is pink.

Data is practically infinite.



- ***>200 BILLION PHOTOS***
- ***>300M PHOTOS UPLOADED/DAY***
- ***3500 PHOTOS EVERY SEC***
- ***ONE IMAGENET EVERY 1:20H***
- ***ONE FLICKR EVERY 4 WEEKS***



Scaling up

DEEPFACE : 4.4 MILLION IMAGES / 4,030 IDENTITIES

RANDOM 108K : 6 MILLION IMAGES / 108,000 IDENTITIES

RANDOM 250K : 10 MILLION IMAGES / 250,000 IDENTITIES

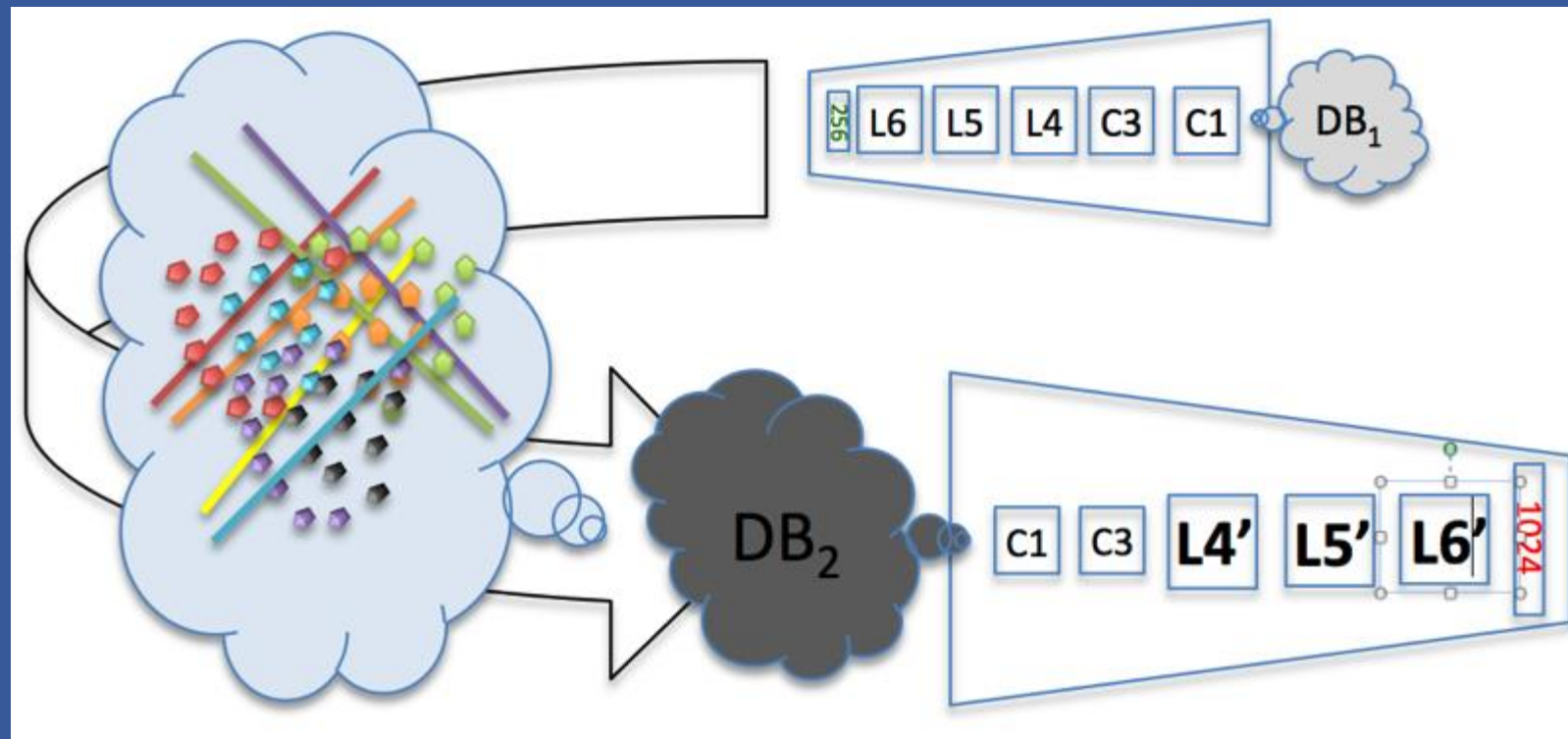
(YES : 250K SOFTMAX)

Training set Dimension	Random 108K			Random 250K			DeepFace
	256	512	1024	256	512	1024	[20]
Verification	97.35	97.62	96.90	96.33	97.10	97.67	97.35

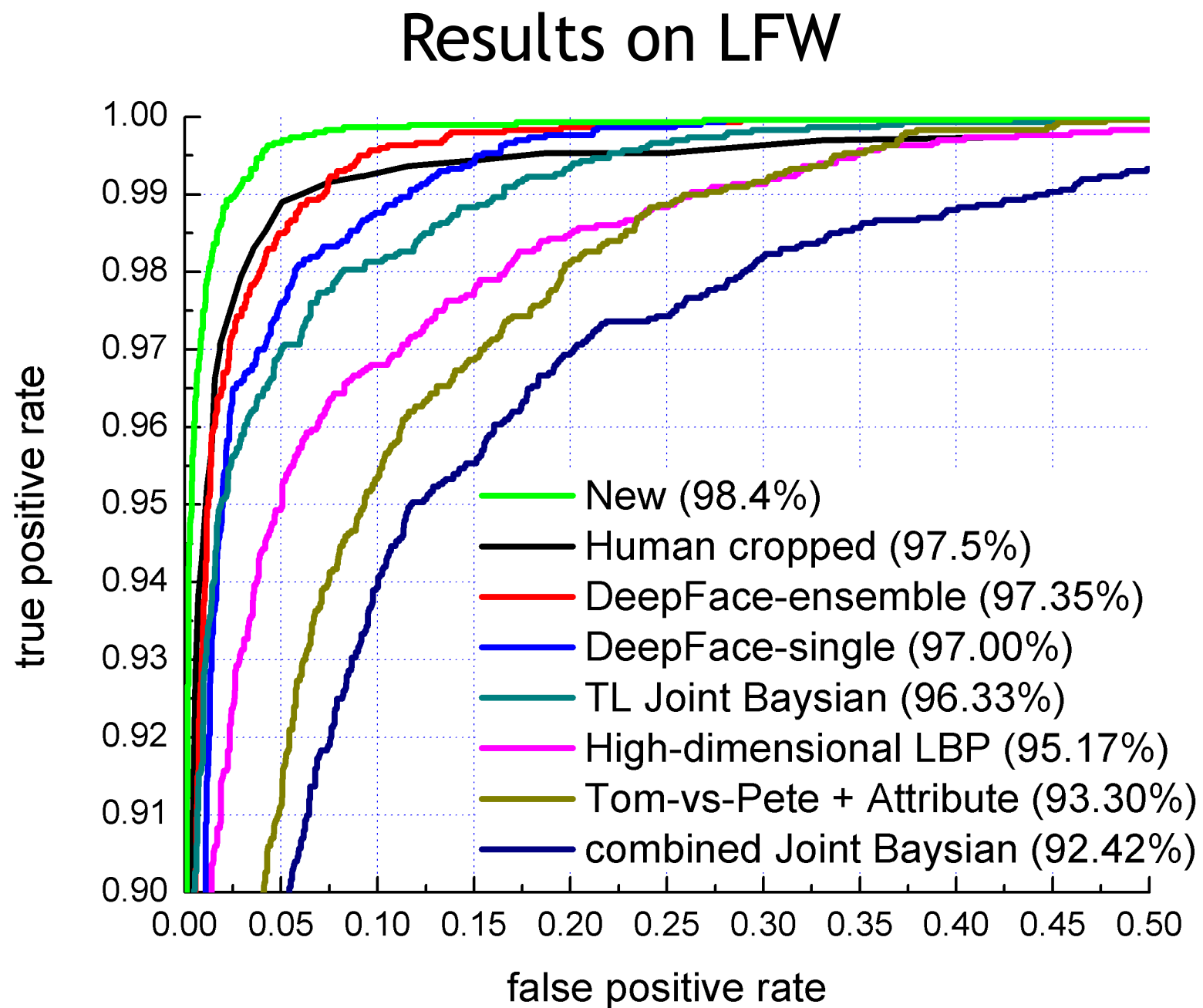
→ SATURATION

Scaling up: Semantic Bootstrapping

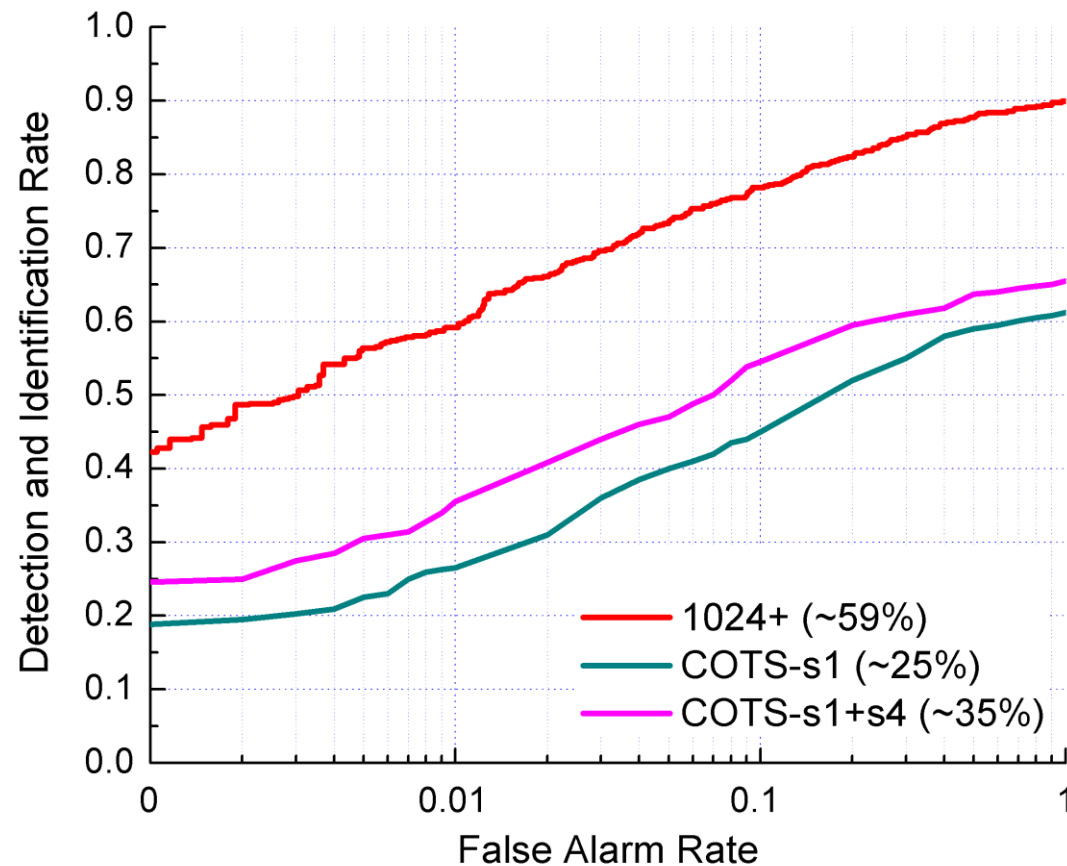
- *0.5B IMAGES → 10M HYPERPLANES*
- *LOOKALIKE HYPERPLANES → DB₂*
- *TRAINING ON DB₂ WITH MORE CAPACITY.*



Second round results



Comparison to NIST's State Of The Art



SECOND-ROUND DEEPFACE

**SAME SYSTEM THAT ACHIEVED
92% RANK-1 ACCURACY ON A TABLE
OF 1.6 MILLION IDENTITIES.
(NIST'S STATE-OF-THE-ART,
CONSTRAINED)**

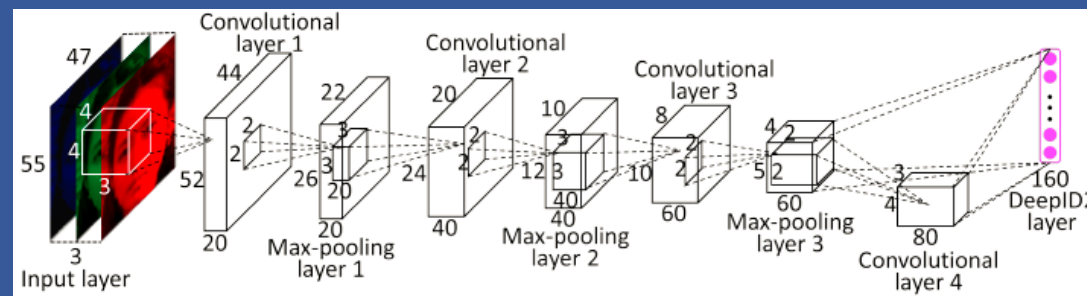
Method	DeepFace [20]	BLS [3]*	COTS- s1 [1]	COTS- s1+s4 [1]	1024+	Fusion
Verification	97.35	93.18	-	-	98.00	98.37
Rank-1	64.9	18.1	56.7	66.5	82.1	82.5
DIR @ 1%	44.5	7.89	25	35	59.2	61.9

DeepFace efficiency (at test)

- For a single 720p image on a single 2.2Ghz Intel CPU core:
 - Face detection: 0.3 sec
 - 2D+ 3D Alignment: 0.05 sec
 - Feed-forward: 0.18 sec
 - Classification: ~0 sec
 - Overall: 0.53 sec
 - Storage: 50%-sparse half-precision floats → ~256 bytes

Additional works

- *Deep learning face representation by joint identification-verification, Sun, Wang, Tang, technical report, arxiv, 6/2014*
- 200 ConvNets from 400 patches ← 2D Aligned (no 3D)
- With Joint Bayesian source / target adaptation
→ 99.15% on the verification (1:1) task.



Additional works

NEW FREE PUBLIC LARGE FACE DATASET FROM SMU:

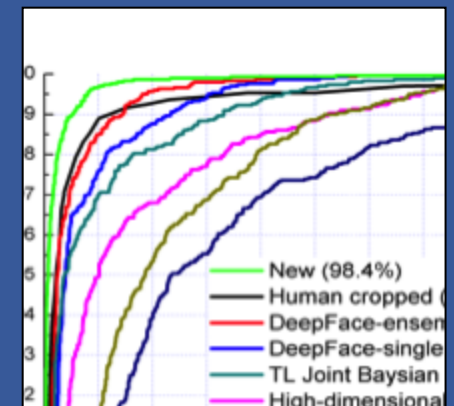
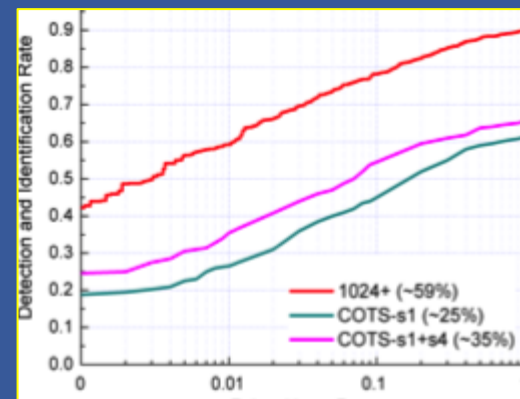
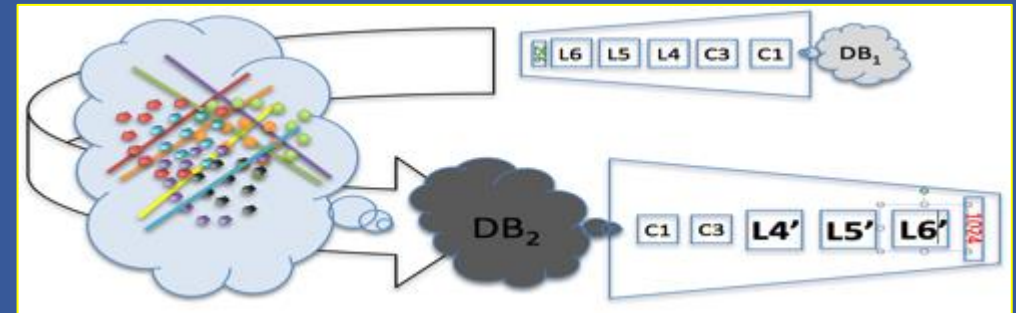
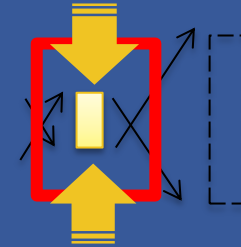
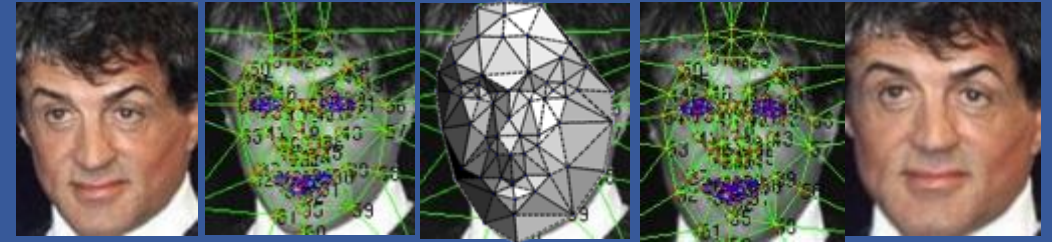
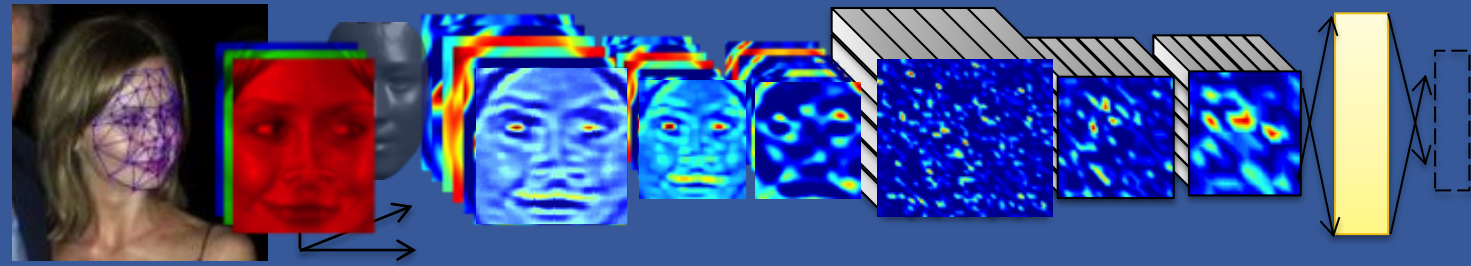
WLFDB : WEAKLY LABELED FACES ON THE WEB

*WANG, DAYONG, HOI, STEVEN C. H., HE, YING, ZHU, JIANKE, MEI, TAO AND LUO, JIEBO,
RETRIEVAL-BASED FACE ANNOTATION BY WEAK LABEL REGULARIZED LOCAL COORDINATE CODING*

714,454 FACIAL IMAGES / 6,025 IDENTITIES

Conclusion:

- Coupling 3D alignment with Large locally-connected networks
- Two-stage 3D alignment system
- Regularization in Transfer Learning
- Scaling up through bootstrapping
- At the brink of human-level performance



Thank you!



¹ Ming Yang



¹ Marc'Aurelio Ranzato



² Lior Wolf

¹ Facebook AI Research

² Tel Aviv University

References:

1. *DeepFace: Closing the Gap to Human-Level Performance in Face Verification*; Taigman, Yang, Ranzato, Wolf
2. *Web-Scale Training for Face Identification*; Taigman, Yang, Ranzato, Wolf
3. *Multi-GPU Training of ConvNets*; Yadan, Adams, Taigman, Ranzato

- Facebook AI Research:
research.facebook.com/ai

