# Center for Brains, Minds & Machines

# Can a biologically-plausible hierarchy effectively replace face detection, alignment, and recognition pipelines?

**by**

Qianli Liao[1], Joel Z Leibo[1], Youssef Mroueh[1], Tomaso Poggio[1]

**Abstract:** The standard approach to unconstrained face recognition in natural photographs is via a detection, alignment, recognition pipeline. While that approach has achieved impressive results, there are several reasons to be dissatisfied with it, among them is its lack of biological plausibility. A recent theory of invariant recognition by feedforward hierarchical networks [1], like HMAX [2, 3], other convolutional networks (e.g., [4]), or possibly the ventral stream, implies an alternative approach to unconstrained face recognition. This approach accomplishes detection and alignment implicitly by storing transformations of training images (called templates) rather than explicitly detecting and aligning faces at test time. Here we propose a particular locality-sensitive hashing based voting scheme which we call "consensus of collisions" and show that it can be used to approximate the full 3-layer hierarchy implied by the theory. The resulting end-to-end system for unconstrained face recognition operates on photographs of faces taken under natural conditions, e.g., Labeled Faces in the Wild (LFW) [5], without aligning or cropping them, as is normally done. It achieves a drastic improvement in the state of the art on this end-to-end task, reaching the same level of performance as the best systems operating on aligned, closely cropped images (no outside training data). It also performs well on two newer datasets, similar to LFW, but more difficult: LFW-jittered (new here) and SUFR-W [6].

arXiv:1311.4082v3 [cs.CV] 26 Mar 2014

---

[1]MIT, McGovern Institute for Brain Research, Center for Brains, Minds and Machines

# 1   Introduction

The challenge of simultaneously ensuring robustness to background clutter and invariance to appearance transformations is pervasive in the design of practical object recognition systems. One approach to mitigating both issues involves subjecting each test image to a detection, alignment, and recognition (DAR) pipeline. The DAR method handles clutter by detecting objects and cropping closely around them so that very little background remains. In a subsequent stage, it attempts to handle transformations by explicit alignment to a standard reference frame—typically effected by rotating and scaling the cropped images to bring certain key features into correspondence. DAR pipelines achieve impressive results in some cases. However, there may be no canonical way to align generic objects, especially if they are non-rigid, thus rendering the approach fundamentally limited.

DAR pipelines are also not plausible models of the brain's recognition system. The brain is not known to have any mechanism of alignment operating on the timescale of visual recognition. Moreover, in computer vision, the problem of unconstrained face recognition is perhaps the best case for the DAR approach since it is relatively simple to align faces to a canonical reference frame using correspondence of internal face features e.g., eyes, nose, and mouth. However, in Biology, face recognition is thought to be "holistic" *i.e.,* not especially driven by the key features predicted by alignment-based strategies [7].

Here we investigate an alternative to the DAR approach based on the idea of using a feedforward hierarchy of stored templates to compute an invariant representation for new input images. This approach is both biologically plausible [8] and theoretically motivated [1]. In order to compare with DAR pipelines, this paper focuses on the problem of unconstrained face recognition. We show here that this approach yields an effective end-to-end system without explicit detection or alignment steps. In particular, we discuss how a system built according to the principles of a recent theory of invariance in hierarchical networks [1] can evade the clutter problem—generally thought to be problematic for feedforward systems [9, 10, 11]. However, use of the system's basic version is limited by the time required to compute full convolutions over space and scale using a large number of filters (50,000 in this case). Thus, the other contribution of this paper is a locality sensitive hashing scheme combined with a voting scheme, that we call *consensus of collisions* (CoC) that approximates the full system implied by the theory. This scheme allows faster computation and scalability of our system.

We argue that the DAR framework guides researchers to focus on each stage of the pipeline in isolation, this is not negative *per se*, but, through the development of specialized datasets researchers into the recognition subproblem come to focus on data that is biased toward the results of particular detection and alignment systems. In particular, Labeled Faces in the Wild-aligned (LFWa) [5, 12], the current gold standard dataset for the recognition step of unconstrained face recognition was filtered by the Viola-Jones face detector [13] and consequently contains almost no faces with any significant rotation in depth. Similarly, alignment was accomplished by a particular commercial system [12] which likely introduces its own subtle biases. Leibo et al. (2014) [6] showed that these biases are severe enough that many recognition systems designed with LFWa in mind do not perform well on a new dataset (SUFR-W), gathered using a very similar protocol to LFW, the primary difference: substituting the, more tolerant of depth-rotation, Zhu & Ramanan face detector [14] for Viola-Jones.

The system proposed here can take as input the full (unaligned, uncropped) images of the LFW dataset. Without performing explicit detection or alignment steps, it achieves a level of performance that compares favorably with the current state of the art systems operating on the aligned and cropped images (87.55%). Note that the present system is solving a much harder problem. The accuracies quoted for the aligned and cropped LFW dataset assume perfect detection and alignment, whereas results using the full images include errors at those (implicit) steps. To further strain the system we also introduced a new "jittered" version of the LFW dataset with additional variation in face position, scale and orientation. The performance of classical systems such as HOG with SVM (Table 4) drops significantly from LFW to LFW-jittered whereas the proposed system is largely unaffected by the additional transformations. We also show that this performance cannot be attributed to overfitting LFW, strong performance is still achieved when the system is trained on SUFR-W and tested on LFW (and vice-versa).

These results demonstrate that the biologically plausible approach is not hopelessly stymied by clutter. Even in the case of unconstrained face recognition, where DAR pipelines could be considered most likely to be effective, this class of biologically plausible hierarchical networks are competitive with the current state of the art end-to-end systems.

# 2   Hierarchical architectures

Hierarchical architectures alternating between selectivity-increasing "tuning" operations and tolerance-increasing "pooling" operations have a long history in computational neuroscience and computer vision [15, 16, 4, 2]. Recently Anselmi et al. (2013) proposed a new theory of these architectures centered around the idea that they compute invariant representations called *signatures*. Two previous face recognition models along the lines suggested by the theory have already been proposed—one was presented as a model of biology [17], the other as a face verification

system [18]. The present proposal incorporates ideas from both while significantly scaling up their scope of operation to the case of unconstrained face recognition without prior detection or alignment.
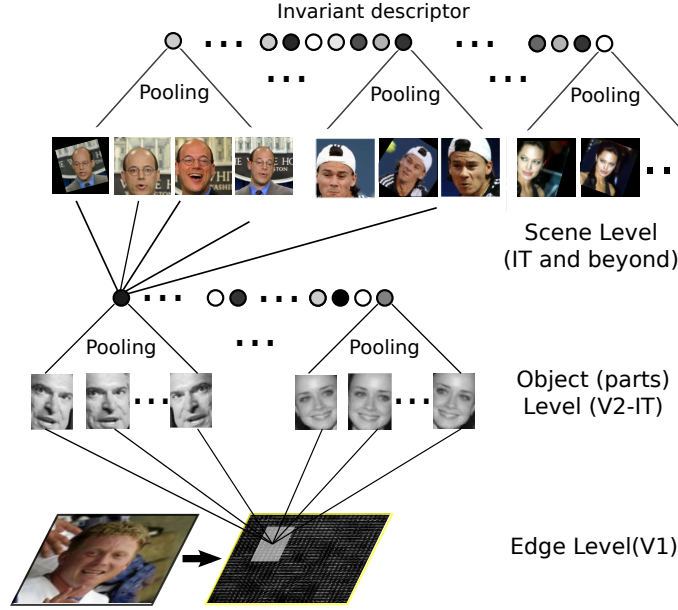


Figure 1: Illustration of the hierarchical architecture

## 2.1 Theoretical motivation

**Background: Invariance and discriminability**

Consider an image $I \in \mathcal{H}$ a Hilbert space. We are interested in recognizing the object depicted by $I$ even when it may have been transformed. Consider a family of transformations $G$ that may (or may not) be a group[2]. The orbit $O_I = \{gI | g \in G\}$ is itself an invariant with respect to the action of $G$. Notice that we use $g$ to refer both to an abstract element of $G$ and to its unitary representation acting on images. For example, if $G$ is the group of in-plane rotations, then the orbit is the set of images generated by rotations of $I$. What if we had instead generated the orbit by transforming $I' = \bar{g}I$? The two orbits are clearly identical, $O_I = O_{I'}$. That is, the set of all rotated images of $I$ is the same as the set of all rotated images of $\bar{g}I$. This motivates a definition. $I$ and $I'$ are considered to be equivalent, written $I \sim I'$, when there exists a $g \in G$ such that $I = gI'$. For example, $I$ and $I'$ would be equivalent if they depict the same object from a different perspective. With this definition, it can be shown that if $G$ is a group then the orbit of any image $I$ under the action of $G$ is unique (see [1]).

Next, let $gI$ be a realization of a random variable. Consider the distribution $P_I$ of images obtained from $I$ under the action of $G$. Anselmi et al. (2013) proved that, for $G$ a group, if two orbits coincide then their associated distributions under $G$ must be identical. This gives the following correspondence between images, orbits, and distributions.

$$I \sim I' \iff O_I = O_{I'} \iff P_I = P_{I'} \tag{1}$$

Thus the distribution $P_I$ is also invariant and unique to each object. The Cramer-Wold theorem [19] suggests a biologically plausible way to characterize such a distribution by its one-dimensional projections. A key result of the theory states (informally) that $P_I$ can be almost uniquely characterized by a set of $K$ one-dimensional distributions $P_{\langle I, t_k \rangle}$ induced by the results of projecting $I$ onto a set of randomly chosen images $t_k$, $k = 1, \ldots, K$ called *templates*. The $P_{\langle I, t_k \rangle}$ can themselves be characterized by their statistical moments, e.g., mean, max, etc. In practice, the number $K$ of projections needed to discriminate a finite number of orbits turns out not to be too large (and Anselmi et al. proved a bound [1]).

Notice however, all this has shown so far is that *if* you had stored (or could compute) $P_{\langle I, t_k \rangle}$, *then* the signature would be invariant and would discriminate between images of different objects. The key fact enabling this approach is that it is possible to store transformations of the templates instead. That is, when $G$ is a group (and $g$ a unitary representation) then

$$\langle gI, t \rangle = \langle I, g^{-1}t \rangle \tag{2}$$

---

[2]For brevity, most of the exposition of the theory given here only applies to compact groups. See [1] for the more general case.
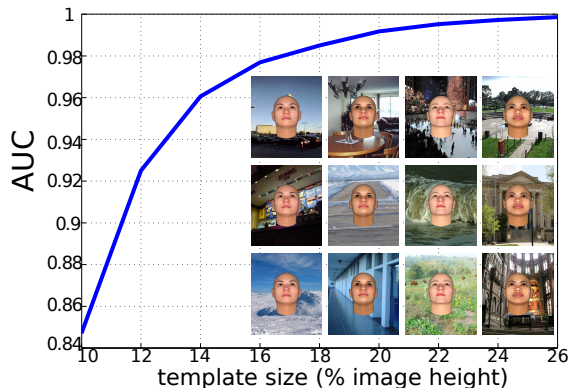
Figure 2: Evaluating clutter tolerance of globally translation invariant layer-2 signatures on a synthetic face verification (same-different matching) task. The only difference between different images of the same person is the background. The dataset is designed to have no other transformation, in order to isolate the clutter problem. The x-axis indicates the template size used, and y axis shows the AUC. Performance increases as templates get increasingly larger. 400 people and 5 background variation per person are used.

This implies that the distributions $P_{\langle gI, t_k \rangle}$ and $P_{\langle I, gt_k \rangle}$ are identical. Thus, either one of them could be used to characterize $P_I$ and uniquely determine the associated orbit. The meaning of the result is as follows. It is not necessary to store all the transformations of $I$, or to have any explicit knowledge of $G$. Rather, storing transformations of the templates is sufficient to construct an invariant for $I$ from a single view.

What happens in the non-group case? It can be shown that under fairly general conditions the signature will be *approximately* invariant (see [1]). However, in this case eq. 2 no longer holds for all $t_k$. There is an additional requirement that $I$ and $t_k$ transform "similarly" (technically this is a condition on the tangent bundles of their respective orbits). Thus approximate invariance for non-group transformations is *class-specific*. In order to compute an invariant signature of $I$ using stored templates in the non-group case, the object depicted in $I$ should be of the same class as the objects depicted in the $t_k$. For example, both $I$ and all the $t_k$ might be faces. [17] conjectured that this class-specificity is the reason the brain's ventral visual pathway separates the processing of faces from other objects [20, 21].

**Background: Architecture**

The architecture consists of a repeated biologically-plausible module that we call an HW-module in honor of Hubel and Wiesel's original proposal for the connectivity of V1 simple and complex cells [15], which was extended as a hypothesis for other ventral stream areas in many computational models e.g., [16, 4, 2]. Here we use the term HW-module to refer to a single "C-unit" and all its afferent "S-units" (note: these units need not correspond to actual cells). A layer of our architecture is a set of HW-modules.

Each HW-module has a single template to which it is "tuned". The "response" $\mu^k(I)$ of the $k$-th HW-module to an image $I$ is given by

$$\mu^k(I) = P_g(\langle I, gt_k \rangle). \tag{3}$$

Where $P_g()$ is a pooling function. The two we will use in the present architecture are the mean and the max (over $g \in G$).

Many different hierarchical architectures consisting of repeated HW-modules are possible. Here we are primarily interested in architectures for which feature complexity as well as invariance increase from early to late layers. Some of these architectures seek to approximately "factorize" image variability into its component transformations. In such architectures, e.g., [17, 18], smaller, edge-like templates are used in early layers while larger templates incorporating information from the entire image are used in higher layers. The early layers discount short-range group transformations while the higher layers compute a representation that is approximately invariant to class-specific transformations. Anselmi et al. (2013) conjectured that these approximately factorizing architectures may improve the sample complexity of multistage learning since invariance in the lower layers could remove the need to align the training images for the higher levels.

### 2.1.1 The clutter problem

A major challenge in the design of any face detection system is the problem of balancing the rate of target acceptance with the rate of false alarms on the background. The hierarchical networks considered here do not have an explicit detection step but they do not escape these issues entirely. A translation-invariant HW-module tuned to a simple template will find high responses at many locations all over any natural image. Whereas a complex template will tend only to be activated when a part of the image is quite similar to it. To illustrate this, we constructed a "pure" test of clutter tolerance using synthetic face images. While these faces clearly do not capture the distribution of natural faces (e.g., none of them have hair), they are convenient for this demonstration. Since the faces themselves are fixed and only the background changes (*i.e.* a model that was not translation invariant and only looked at the center of the image would always perform perfectly) the only way to fail is due to spurious activations on the background. Increasing the size of the templates steadily improves the performance on this pure test of clutter tolerance (fig. 2). Similar results with several other translation invariant architectures have been reported before, e.g., SIFT variants [22, 6] so this is likely to be a general issue with translation invariance and not a quirk of our particular system or dataset.

---

**Algorithm 1** Consensus of Collisions

---
**Input:** Test images $I$, templates $\{t_j\}_{j=1\ldots n}$, number of consensus $N$.
**Notations:** $W$ : Oversampled windows, $C$ : Candidate Windows , $S$ : Consensus windows, $\langle,\rangle$ : normalized dot product.
**Output:** $R$ : the Response of templates for all images.
**Code:**
**Training:**
Initialize hash function h
Compute hash codes of $H_{t_j} = h(t_j) \quad \forall j$
**Testing:**
**for** $i = 1$ **to #images do**
  W $\leftarrow$ Dense oversampling of windows from image $I_i$
  Compute hash codes $H_w = h(w), \quad \forall w \in W$.
  **for** $j = 1$ **to m do**
    $C_{t_j} \leftarrow \{w \in W, |h(w) = h(t_j)\}$
  **end for**
  $S \leftarrow$ Select $N$ most frequently appearing windows out of $\cup_j C_{t_j}$
  **for** $j = 1$ **to n do**
    Set $P$ to be a empty Vector
    **for** all windows $w$ in $S$ **do**
      $P(k) = \langle w, t_j \rangle$
      $k = k + 1$
    **end for**
    $R(i,j) =$`max`$(P)$
  **end for**
**end for**

---

## 3 Architecture and Approximations

Given a pair of images $(x_a, x_b)$ the task is to verify whether they depict the same person or not. To test our HW-architecture we run it on both images and compare them by the angle between their signatures (top-level representations). That is, we take the normalized dot product $\langle \mu(x_a), \mu(x_b) \rangle$ if it exceeds a threshold $\tau$, our method outputs that $x_a$ and $x_b$ have the same identity otherwise they are different. We use the training set to pick the optimal $\tau$.

In this section we detail our full architecture given in figure 1. Since the architecture consists of a hierarchy of HW-modules, it can be thought of as a succession of simple and complex cells performing two main operations tuning (projection on a template) and pooling. The final output is the signature $\mu(I)$, a vector of top-level HW-module responses, each tuned to the identity of a face, invariantly to affine (group) transformations and approximately invariant to class-specific transformations.

1. **First Layer:** The proposed system uses closely cropped face images for training (but not for testing). For each of the $n$ closely cropped face templates, we compute low-level features at each position/scale. These

could be either HOG [23] features, or in some cases: HOG and LBP [24] features were combined. In those cases, the second layer was computed separately for each feature type and the results fused by concatenation before computing the third layer.. For training we extract $n$ HoG templates from *closely cropped face images*. We call those templates the second layer training templates.

2. **Second Layer:** We extract for each test image a dense overlapping set of $m$ windows. We convolve the second layer training templates with all the windows, and then apply max pooling over all scales and locations. For each template, we pool the responses over all windows. Finally, for templates generated from a group of transformations (in-plane rotation in this case) our system also does max pooling over that group[3].

3. **Third Layer:** For each training image, run the architecture until the second layer. Store the responses up to the second layer and use them as the third-layer training templates. For a test image, compute the dot product of the output of the second layer, with the stored third layer training templates. Note that the third layer training templates are indexed by the identity of the person. Thus, as in [17, 18], each third-level HW-module pools over a set of templates depicting the same person.
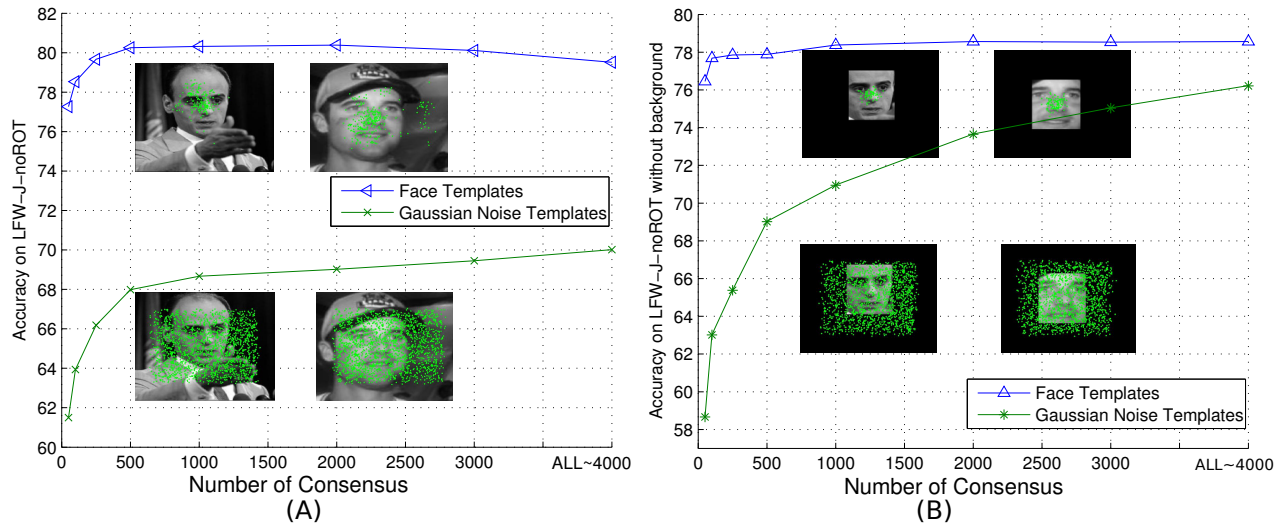


Figure 3: Evaluating the effect of "Consensus of Collisions" and background tolerance of class-specific templates. A set of experiments are run on two LFW-jittered-noRot datasets (for simplicity, no in-plane rotation jittering)— with or without background (figure A and B respectively). Each experiment used 5,000 templates of either faces or gaussian noise. The green dots in the images are examples of the maximum response locations of the templates of each curve (edges are blank because only valid convolution locations are considered). The last data point of each curve shows the performance when no CoC is performed.

## 3.1 Approximation by locality sensitive hashing and low rank approximation

Layer two is the computational bottleneck. It is computationally / memory-access expensive to compute normalized dot products. Consider a set of windows $(w_1 \ldots w_m)$ representing an image by its sub-patches at a range of positions and scales. The problem is to compute the normalised dot product $\langle w_i, t_j \rangle$ between each window $w_i$ and each template $t_j$ in a large collection.

Inspired by the impressive results of Dean et al. [25], we experimented with approximating our architecture by locality sensitive hashing (LSH). Specifically, we propose a simple LSH scheme called "Consensus of Collisions" (CoC) (Algorithm 1) that evades the need for intensive memory accessing for dot product computations and reduces the computational load. Let $h$ denote a randomised hash function. CoC proceeds in three steps in order to prune windows where it is unlikely to find a face:

1. **Hashing:** Produce hash tables of templates $(h(t_1) \ldots h(t_n))$ and windows $(h(w_1) \ldots h(w_m))$.

2. **Collect candidate "match" windows for each template:** For each template $t_j$, find a set of candidate windows that have exactly the same hash code. Let $C_{t_j}$ be that set.

---

[3]Via eq. 2, for group transformations, transforming the test image or the template is equivalent.

3. **Voting Scheme:** For each window $w_i$ compute the frequency of being chosen across templates, $f_{w_i} = \frac{1}{n} \sum_{j=1}^{n} 1_{w_i \in C_{t_j}}$ . Pick the $N$ "most popular" windows across templates i.e the $N$ windows having the highest frequencies. Let $S$ be the set of the $N$ most popular windows . We call $N$ *the number of consensus*, and $S$ *the consensus set.*

This locality sensitive hashing and voting scheme allow us to efficiently find a set of $N$ windows with high responses across templates. Note that $N \ll m$, so the speedup due to CoC is quite large.

The final step is to pool across the windows of the consensus set $S$. For each template $t_j$, find the window in $S$ with the highest response i.e., $\max_{w \in S} \frac{\langle t_j, w \rangle}{||t_j||||w||}$. We explored two options for computing the normalised dot products at this stage: 1. Exact computation, or 2., The PCA approximation described next.

**PCA approximation**

Note that the matrix $T$, consisting of all the templates (as column or row vectors), is a low rank matrix. We can perform PCA on $T$ and keep the $k$ largest eigenvectors. By projecting the templates and the windows to the $k$ dimensional space defined by those eigenvectors we can perform faster dot products in the reduced dimensional space. This procedure is similar to the one employed by [26].

# 4    Experimental Evaluation

In this section, we demonstrate the performance and properties of our approach through a set of experiments on LFW, SUFR-W and a series of difficult LFW-jittered datasets we created.

## 4.1    Datasets

We used the following five datasets throughout the paper:(1) LFW-original: The full (250x250), non-cropped, non-aligned version of LFW dataset (Figure. 4). (2) SUFR-W: Unconstrained face recognition dataset collected by [6] with similar protocol to LFW but with a more advanced detector [14]. (3) LFW-Jittered (LFW-J): The LFW-Jittered dataset was created by randomly translating, scaling and in-plane rotating LFW original images. Translation range: -40 to 40 pixels, scaling range: 1 to 1.5, in-plane rotation range: -20 to 20 degrees. See the supplementary information for more details on the datasets and example images.

## 4.2    Experiments

**Full model:** We tested the full model described in Figure 1 on the LFW-original dataset without using any hashing approximation (first two rows of Table 1).

**Approximated model:** Despite the encouraging accuracy of the direct implementation, it is not fast enough for practical purposes. Thus, we further explored the effect of hashing and several different choices of PCA and feature type (the rest of table 1). The system achieved 82.53% performance on the original LFW data with a system that runs at nearly 2 frames per second. Note that hashing becomes indispensable when the windows number becomes large. We tested a model with a large number of scales—the pyramid contains 31-scales from 125x125 (50%) to 500x500 (200%), generating about 30,000 windows at test time. In this case, any model without hashing is memory-intractable (requires >25 GB per thread), and hashing alone gives a 60x speedup.

**Large class-specific templates:** As suggested by fig. 2, large templates mitigate the clutter problem of translation-invariant HW-modules. To experimentally demonstrate this property, a set of experiments were run on LFW-jittered-noRot and LFW-jittered-noRot-noBG. For each experiment, we used 5,000 templates of either faces or gaussian noise. Performances are plotted as a function of the number of consensus $N$ kept (Figure 3). The performance of gaussian noise templates is surprisingly high when no clutter is present, but much lower otherwise. With face templates, the model tolerates significant clutter and is highly selective to faces, in which case hashing is able to reduce the number of windows computed by 90% without lowering the performance.

**Robustness across datasets:** Despite the recent close-to-human performance reported in LFW-a, (e.g., [27]), it has been argued by [6] that good LFW performance often does not transfer to SUFR-W, indicating that the community may be somewhat overfitting LFW. To demonstrate the robustness of our model, we tested its performance across LFW and SUFR-W. The approach is to train on one dataset and test on the other. Since our model has two stages of trainable templates, we tried either partial or full training using the other dataset's training set. The performance is shown in Table 2. The findings are: 1. training the first layer on either dataset gives similar

| Feature | #EigenVec. | #Consensus | Mem. | Time. | Acc. | Speedup |
|---------|-----------|-----------|------|-------|------|---------|
| LBP | No PCA | No Hashing | 6.8 GB | 44.67 | 83.67 | 1x |
| LBP | No PCA | 500 | 2.9 GB | 4.38 | 82.40 | 10.2x |
| LBP | 250 | No Hashing | 2.1 GB | 4.13 | 82.18 | 10.8x |
| LBP | 1200 | 1500 | 2.0GB | 2.51 | 83.17 | 17.7x |
| HOG | 1200 | 1500 | 1.8GB | 1.69 | **84.73** | 26.4x |
| LBP | 250 | 500 | 1.5 GB | 0.89 | 81.18 | 50.2x |
| HOG | 250 | 500 | 1.1 GB | **0.54** | 82.53 | **82.7x** |

Table 1: The performance (Acc.) evaluated on the unaligned LFW dataset. There were about 50,000 face templates in our model. The testing image (250x250) was scaled to form a 12-scale pyramid with size ranging from 288x288 to 150x150. The dimensionalities of the LBP [24] and HOG [23] features were 7540 and 4030 respectively. Every experiment was run on a single 2010 machine with an 8-core processor and 36GB RAM. No GPU was used. Multi-threading was employed to use the CPU as much as possible. The time refers to the average time spent on a single frame (13233 frames are tested in total). The "Mem" refers to the average memory requirement of each thread. The code was written in Matlab and still has optimization potential. A GPU implementation was found to be inefficient due to the GPU's small on-board memory. The number of CPU cores is the bottleneck. Further speedup is expected if more of them are available. Hash code length: 24, Hash table number: 20. Note: The training process just performs PCA and stores templates—it only takes about 5 to 10 minutes (contrast with most deep learning approaches (e. g., [4]).

| 2nd Layer | 3rd Layer | Test on | Acc. |
|-----------|-----------|---------|------|
| SUFR-W | SUFR-W | SUFR-W | 80.30±0.89% |
| LFW | SUFR-W | SUFR-W | 79.60±1.41% |
| LFW | LFW | SUFR-W | 76.08±1.36% |
| LFW | LFW | LFW | 84.33±1.79% |
| SUFR-W | LFW | LFW | 83.87±1.38% |
| SUFR-W | SUFR-W | LFW | 84.55±1.43% |

Table 2: Dataset bias: training 2nd (and 3rd) layer(s) with data from another dataset. SUFR-W works better as the third layer because it has more people (400), while LFW only has about 150 people with more than 10 images. Model description: HOG features, #eigen vector 1200, #consensus 1500, Hash code length 28, Hash table number 20.

performance. 2. LFW is significantly easier than SUFR-W.

**The LFW-original state-of-the-art:** In Table 3, we demonstrate our system's performance on LFW, comparing to the "no outside data used, unaligned" category. Note however, since our system uses the identities of the faces in the training set, it does not exactly conform to the recommendations of [5] for "restricted" training with LFW which requires only same-different pairs to be used. Our procedure follows the "restricted" protocol at test time, but not for training. We intend to model what the brain could learn via unsupervised mechanisms—these template orbits could have been learned in an unsupervised fashion by observing faces transform and pooling over temporally adjacent frames. The same setting is addressed in [17, 18, 1].

**Jittering Invariance:** Human vision is invariant to small shifts in object position. Motivated by the strong performance on LFW-original, we tested the same model on LFW-jittered data (Table 4), which we expected to be very difficult for conventional computer vision methods. Our system achieved almost the same performance on this dataset. In contrast, the baseline model (HOG) drops by almost 20%. Intriguingly, our training templates were only rotated between -12 to 12 degrees, but the model handles -20 to 20 jittering in LFW-J without any problem. This is likely due to preexisting angle variation in the training templates (they were not aligned). As in [18], Non-uniform and ultra-sparse sampling of the orbit are sufficient for good performance [18].

# 5 Conclusion

What then is the answer to this paper's title question? Can the proposed network effectively replace DAR pipelines for unconstrained face recognition? The results are promising in that direction. Our system achieves one of the strongest reported accuracies in LFW's "unaligned & no outside data used" category (fig. 3). It also performs well on two more difficult datasets: SUFR-W and a significantly jittered (misaligned) version of LFW (example images

| LFW no outside data used | | | | | |
| --- | --- | --- | --- | --- | --- |
| Aligned | | Unaligned | | | |
| **Model** | **Acc.** | **Model** | **Acc.** | **Model (translation-invar.)** | **Acc.** |
| Wolf et al. [28] | 78.47% | Nowak et al. [29] | 72.45% | SIFT-BoW+SVM(Baseline) | 57.73±2.53% |
| V1-like/MKL [30] | 79.35% | Sanderson et al. [31] | 72.95% | Our Model(HOG) | 84.73±1.82% |
| APEM (fusion) [32] | 84.08% | MRF-MLBP [33] | 79.08% | Our Model(HOG+LBP) | **86.15±1.50%** |
| Simonyan et al. [34] | 87.47% | APEM (fusion) [32] | 81.70% | Our Model(HOG+LBP)+SVM | **87.55±1.41%** |

Table 3: Our model (87.55±1.41%) significantly outperforms state-of-the-art: APEM (81.70±1.78%) in the LFW "unaligned & no outside data used" category. The last column shows models that are translation-invariant. For the last model, we simply replace the final cosine distance classifier with a RBF-SVM, and we used the difference of the testing pair's third layer signatures as the input to the SVM.

| Model | LFW | LFW-J |
| --- | --- | --- |
| HOG+SVM (baseline) | 74.45/67.32% | 55.28% |
| Our Model (HOG) | 84.73% | 84.62% |
| Our Model (HOG+LBP) | 86.15% | 86.02% |
| Our Model (HOG+LBP) + SVM | 87.55% | 87.45% |

Table 4: The performance on LFW-original (unaligned) and LFW-J (jittered) datasets. For LFW-J: translation range: -40 to 40 pixels, scaling range: 1 to 1.5, in-plane rotation range -20 to 20 degree. We used *exactly* the same model for LFW and LFW-J. The HOG baseline, 74.45% uses the closely cropped performance and 67.32% is non-cropped performance. With jittering, one cannot crop the image. Either way, HOG performance drops dramatically.

in fig. 4).

Another encouraging result is the proposed system's robustness to the choice of dataset from which to obtain training templates (fig. 2). This suggests that its strong performance is unlikely to be attributable to overfitting. The finding that it is better to train on SUFR-W and test on LFW than vice-versa is likely due to the presence of non-frontal faces in the former dataset but not the latter.

The proposed hierarchical approach has applications in computational neuroscience: 1. as a proof of principle that a feedforward network is able to do unconstrained face recognition, 2., as a starting point for developing models of the ventral stream and its face-specific branch[4] and 3., as a justification for the use of large templates to mitigate the problem of clutter—which is particularly interesting since large templates are known to provide an explanation for holistic face effects [36] like the composite face effect [37].

We developed our system in order to address the question of whether a unified hierarchy could perform competitively with the dominant DAR approach. However, it is also interesting to consider the possibility of synergies between the two. Nothing about our method precludes its inclusion *as* the recognition (or recognition + alignment) module within a standard DAR pipeline. Such a hybrid system may be able to recover from errors in previous pipeline stages. For example, an invariant recognition module may be able to rescue a positive result despite an error in the alignment stage. We think this is a promising avenue for future investigation.

# Acknowledgments

# References

[1] Anselmi, F., Leibo, J.Z., Rosasco, L., Mutch, J., Tacchetti, A., Poggio, T.: Unsupervised learning of invariant representations in hierarchical architectures. arXiv preprint arXiv:1311.4158 (2013)

[2] Riesenhuber, M., Poggio, T.: Hierarchical models of object recognition in cortex. Nature Neuroscience **2**(11) (November 1999) 1019–1025

---

[4]The feedforward network of face-specific patches of visual cortex described by [21, 35].

(A) LFW-a    (B) LFW-original    (C) LFW-Jittered

Figure 4: (A) LFW-a dataset—The majority of studies on LFW actually use the more finely aligned dataset LFW-a and crop very close as shown here [6]. (B) Original LFW images. (C) Our LFW-jittered dataset.

[3] Serre, T., Wolf, L., Bileschi, S., Riesenhuber, M., Poggio, T.: Robust Object Recognition with Cortex-Like Mechanisms. IEEE Transactions on Pattern Analysis and Machine Intelligence **29**(3) (2007) 411–426

[4] LeCun, Y., Bengio, Y.: Convolutional networks for images, speech, and time series. The handbook of brain theory and neural networks (1995) 255–258

[5] Huang, G.B., Mattar, M., Berg, T., Learned-Miller, E.: Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In: Workshop on faces in real-life images: Detection, alignment and recognition (ECCV), Marseille, Fr (2008)

[6] Leibo, J.Z., Liao, Q., Poggio, T.: Subtasks of Unconstrained Face Recognition. In: International Joint Conference on Computer Vision, Imaging and Computer Graphics, VISIGRAPP, Lisbon, Portugal (2014)

[7] Tanaka, J., Farah, M.: Parts and wholes in face recognition. The Quarterly Journal of Experimental Psychology **46**(2) (1993) 225–245

[8] DiCarlo, J.J., Zoccolan, D., Rust, N.C.: How does the brain solve visual object recognition? Neuron **73**(3) (2012) 415–434

[9] Treisman, A.: Features and objects: The fourteenth Bartlett memorial lecture. The Quarterly Journal of Experimental Psychology **40**(2) (1988) 201–237

[10] Itti, L., Koch, C.: Computational modelling of visual attention. Nature Reviews Neuroscience **2**(3) (2001) 194–203

[11] Chikkerur, S.S., Serre, T., Tan, C., Poggio, T.: What and where: A Bayesian inference theory of attention. Vision Research (May 2010)

[12] Taigman, Y., Wolf, L., Hassner, T.: Multiple One-Shots for Utilizing Class Label Information. In: British Machine Vision Conference. (2009) 1–12

[13] Viola, P., Jones, M.J.: Robust real-time face detection. International journal of computer vision **57**(2) (2004) 137–154

[14] Zhu, X., Ramanan, D.: Face detection, pose estimation, and landmark localization in the wild. In: IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), Providence, RI (2012) 2879–2886

[15] Hubel, D., Wiesel, T.: Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. The Journal of Physiology **160**(1) (1962) 106

[16] Fukushima, K.: Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. Biological Cybernetics **36**(4) (April 1980) 193–202

[17] Leibo, J.Z., Mutch, J., Poggio, T.: Why The Brain Separates Face Recognition From Object Recognition. In: Advances in Neural Information Processing Systems (NIPS), Granada, Spain (2011)

[18] Liao, Q., Leibo, J.Z., Poggio, T.: Learning invariant representations and applications to face verification. In: Advances in Neural Information Processing Systems (NIPS), Lake Tahoe, CA (2013)

[19] Cramér, H., Wold, H.: Some theorems on distribution functions. Journal of the London Mathematical Society **1**(4) (1936) 290–294

[20] Kanwisher, N., McDermott, J., Chun, M.: The fusiform face area: a module in human extrastriate cortex specialized for face perception. The Journal of Neuroscience **17**(11) (1997) 4302

[21] Tsao, D., Freiwald, W.A., Tootell, R., Livingstone, M.: A cortical region consisting entirely of face-selective cells. Science **311**(5761) (2006) 670

[22] Ruiz-del Solar, J., Verschae, R., Correa, M.: Recognition of faces in unconstrained environments: a comparative study. EURASIP Journal on Advances in Signal Processing **2009** (2009)

[23] Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. IEEE International Conference on Computer Vision and Pattern Recognition (CVPR) **1**(886-893) (2005)

[24] Ojala, T., Pietkainen, M., Maenpaa, T.: Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. IEEE Transactions on Pattern Analysis and Machine Intelligence **24**(7) (2002) 971–987

[25] Dean, T., Ruzon, M.A., Segal, M., Shlens, J., Vijayanarasimhan, S., Yagnik, J.: Fast, Accurate Detection of 100,000 Object Classes on a Single Machine, url = , year = 2013. . . . IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2013

[26] Chikkerur, S., Poggio, T.: Approximations in the hmax model. MIT-CSAIL-TR-2011-021, CBCL-298 (2011)

[27] Chen, D., Cao, X., Wen, F., Sun, J.: Blessing of Dimensionality: High-dimensional Feature and Its Efficient Compression for Face Verification. In: IEEE International Conference on Computer Vision and Pattern Recognition (CVPR). (2013)

[28] Wolf, L., Hassner, T., Taigman, Y., et al.: Descriptor based methods in the wild. In: Workshop on Faces in'Real-Life'Images: Detection, Alignment, and Recognition. (2008)

[29] Nowak, E., Jurie, F.: Learning visual similarity measures for comparing never seen objects. In: Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on, IEEE (2007) 1–8

[30] Pinto, N., DiCarlo, J.J., Cox, D.D.: How far can you get with a modern face recognition test set using only simple features? In: Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on, IEEE (2009) 2591–2598

[31] Sanderson, C., Lovell, B.C.: Multi-region probabilistic histograms for robust and scalable identity inference. In: Advances in Biometrics. Springer (2009) 199–208

[32] Cui, Z., Li, W., Xu, D., Shan, S., Chen, X.: Fusing robust face region descriptors via multiple metric learning for face recognition in the wild. In: Computer Vision and Pattern Recognition (CVPR). (2013)

[33] Arashloo, S.R., Kittler, J.: Efficient processing of mrfs for unconstrained-pose face recognition. In: Biometrics: Theory, Applications and Systems. (2013)

[34] Karen Simonyan, Omkar M. Parkhi, A.V., Zisserman, A.: Fisher vector faces in the wild. In: British Machine Vision Conference (BMVC). (2013)

[35] Freiwald, W.A., Tsao, D.: Functional Compartmentalization and Viewpoint Generalization Within the Macaque Face-Processing System. Science **330**(6005) (2010) 845

[36] Tan, C., Poggio, T.: Faces as a "Model Category" for Visual Object Recognition. MIT-CSAIL-TR-2013-004, CBCL-311 (2013)

[37] Young, A.W., Hellawell, D., Hay, D.C.: Configurational information in face perception. Perception **16**(6) (1987) 747–759