

# Flowing ConvNets for Human Pose Estimation in Videos

Tomas Pfister  
Dept. of Engineering Science  
University of Oxford  
tp@robots.ox.ac.uk

James Charles  
School of Computing  
University of Leeds  
j.charles@leeds.ac.uk

Andrew Zisserman  
Dept. of Engineering Science  
University of Oxford  
az@robots.ox.ac.uk

## Abstract

*The objective of this work is human pose estimation in videos, where multiple frames are available. We investigate a ConvNet architecture that is able to benefit from temporal context by combining information across the multiple frames using optical flow.*

*To this end we propose a network architecture with the following novelties: (i) a deeper network than previously investigated for regressing heatmaps; (ii) spatial fusion layers that learn an implicit spatial model; (iii) optical flow is used to align heatmap predictions from neighbouring frames; and (iv) a final parametric pooling layer which learns to combine the aligned heatmaps into a pooled confidence map.*

*We show that this architecture outperforms a number of others, including one that uses optical flow solely at the input layers, one that regresses joint coordinates directly, and one that predicts heatmaps without spatial fusion.*

*The new architecture outperforms the state of the art by a large margin on three video pose estimation datasets, including the very challenging Poses in the Wild dataset, and outperforms other deep methods that don't use a graphical model on the single-image FLIC benchmark (and also [5, 35] in the high precision region).*

## 1. Introduction

Despite a long history of research, human pose estimation in videos remains a very challenging task in computer vision. Compared to still image pose estimation, the temporal component of videos provides an additional (and important) cue for recognition, as strong dependencies of pose positions exist between temporally close video frames.

In this work we propose a new approach for using optical flow for part localisation in deep Convolutional Networks (ConvNets), and demonstrate its performance for human pose estimation in videos. The key insight is that, since for localisation the prediction targets are positions in the image space (e.g.  $(x, y)$  coordinates of joints), one can use

dense optical flow vectors to *warp predicted positions* onto a target image. In particular, we show that when regressing a *heatmap* of positions (in our application for human joints), the heatmaps from neighbouring frames can be warped and aligned using optical flow, effectively propagating position confidences temporally, as illustrated in Fig 1.

We also propose a deeper architecture that has additional convolutional layers beyond the initial heatmaps to enable learning an *implicit* spatial model of human layout. These layers are able to learn dependencies between human body parts. We show that these ‘spatial fusion’ layers remove pose estimation failures that are kinematically impossible.

**Related work.** Traditional methods for pose estimation have often used pictorial structure models [2, 8, 10, 27, 39], which optimise a configuration of parts as a function of local image evidence for a part, and a prior for the relative positions of parts in the human kinematic chain. An alternative approach uses poselets [1, 13]. More recent work has tackled pose estimation holistically: initially with Random Forests on depth data [12, 29, 31, 34] and RGB [3, 4, 24], and most recently with convolutional neural networks.

The power of ConvNets has been demonstrated in a wide variety of vision tasks – object classification and detection [11, 21, 28, 40], face recognition [32], text recognition [15, 16], video action recognition [20, 30] and many more [7, 22, 25].

For pose estimation, there were early examples of using ConvNets for pose comparisons [33]. More recently, [37] used an AlexNet-like ConvNet to directly regress *joint coordinates*, with a cascade of ConvNet regressors to improve accuracy over a single pose regressor network. Chen and Yuille [5] combine a parts-based model with ConvNets (by using a ConvNet to learn conditional probabilities for the presence of parts and their spatial relationship with image patches). In a series of papers, Tompson, Jain *et al.* developed ConvNet architectures to directly regress *heatmaps* for each joint, with subsequent layers to add an Markov Random Field (MRF)-based spatial model [17, 36], and a pose refinement model [35] (based on a Siamese network with shared weights) upon a rougher pose estimator ConvNet.

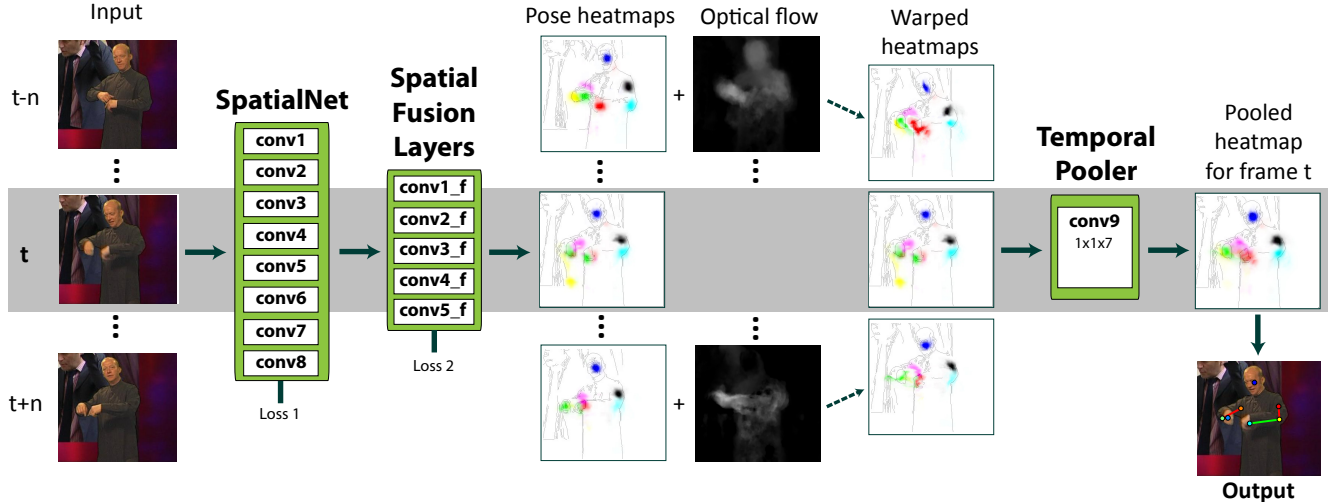


Figure 1. **Deep expert pooling architecture for pose estimation.** The network takes as an input all RGB frames within a  $n$ -frame neighbourhood of the current frame  $t$ . The fully convolutional network (consisting of a heatmap net with an implicit spatial model) predicts a confidence heatmap for each body joint in these frames (shown here with a different colour per joint). These heatmaps are then temporally *warped* to the current frame  $t$  using optical flow. The warped heatmaps (from multiple frames) are then *pooled* with another convolutional layer (the temporal pooler), which learns how to weigh the warped heatmaps from nearby frames. The final body joints are selected as the maximum of the pooled heatmap (illustrated here with a skeleton overlaid on top of the person).

Temporal information in videos was initially used with ConvNets for action recognition [30], where optical flow was used as an input *motion feature* to the network. Following this work, [18, 24] investigated the use of temporal information for pose estimation in a similar manner, by inputting flow or RGB from multiple nearby frames into the network, and predicting joint positions in the current frame.

However, pose estimation differs from action recognition in a key respect which warrants a different approach to using optical flow: in action recognition the prediction target is a class label, whereas in pose estimation the target is a set of  $(x, y)$  positions onto the image. Since the targets are positions in the image space, one can use dense optical flow vectors not only as an input feature but also to *warp predicted positions* in the image, as done in [4] for random forest estimators. To this end, our work explicitly predicts joint positions for *all* neighbouring frames, and temporally aligns them to frame  $t$  by warping them backwards or forwards in time using tracks from dense optical flow. This effectively reinforces the confidence in frame  $t$  with a strong set of ‘expert opinions’ (with corresponding confidences) from neighbouring frames, from which joint positions can be more precisely estimated. Unlike [4] who average the expert opinions, we learn the expert pooling weights with backpropagation in an end-to-end ConvNet.

Our ConvNet outperforms the state of the art on three challenging video pose estimation datasets (BBC Pose, ChaLearn and Poses in the Wild) – the heatmap regressor alone surpasses the state of the art on these datasets, and the pooling from neighbouring frames using optical flow gives

a further significant boost. We have released the models and code at <http://www.robots.ox.ac.uk/~vgg>.

## 2. Temporal Pose Estimation Networks

Fig 1 shows an overview of the ConvNet architecture. Given a set of input frames within a temporal neighbourhood of  $n$  frames from a frame  $t$ , a spatial ConvNet regresses joint confidence maps (‘heatmaps’) for each input frame separately. These heatmaps are then individually *warped* to frame  $t$  using dense optical flow. The warped heatmaps (which are effectively ‘expert opinions’ about joint positions from the past and future) are then pooled into a single heatmap for each joint, from which the pose is estimated as the maximum.

We next discuss the architecture of the ConvNets in detail. This is followed by a description of how optical flow is used to warp and pool the output from the Spatial ConvNet.

### 2.1. Spatial ConvNet

The network is trained to regress the location of the human joint positions. However, instead of regressing the joint  $(x, y)$  positions directly [24, 37], we regress a *heatmap* of the joint positions, separately for each joint in an input image. This heatmap (the output of last convolutional layer, conv8) is a fixed-size  $i \times j \times k$ -dimensional cube (here  $64 \times 64 \times 7$  for  $k = 7$  upper-body joints). At training time, the ground truth label are heatmaps synthesised for each joint separately by placing a Gaussian with fixed variance at the ground truth joint position (see Fig 2). We then

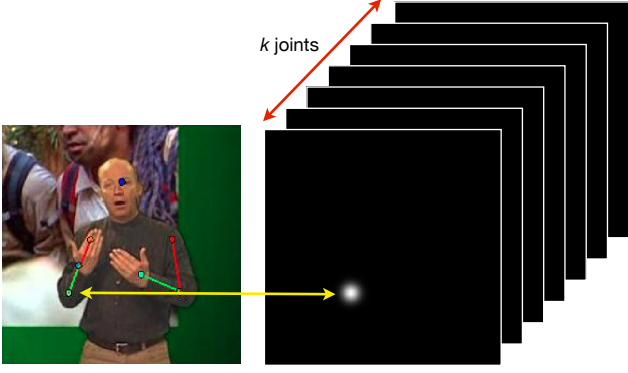


Figure 2. **Regression target for learning the Spatial ConvNet.** The learning target for the convolutional network is (for each of  $k$  joints) a heatmap with a synthesised Gaussian with a fixed variance centred at the ground truth joint position. The loss is  $l_2$  between this target and the output of the last convolutional layer.

use an  $l_2$  loss, which penalises the squared pixel-wise differences between the predicted heatmap and the synthesised ground truth heatmap.

We denote the training example as  $(\mathbf{X}, \mathbf{y})$ , where  $\mathbf{y}$  stands for the coordinates of the  $k$  joints in the image  $\mathbf{X}$ . Given training data  $N = \{\mathbf{X}, \mathbf{y}\}$  and the ConvNet regressor  $\phi$  (output from conv8), the training objective becomes the task of estimating the network weights  $\lambda$ :

$$\arg \min_{\lambda} \sum_{(\mathbf{X}, \mathbf{y}) \in N} \sum_{i,j,k} \|G_{i,j,k}(\mathbf{y}_k) - \phi_{i,j,k}(\mathbf{X}, \lambda)\|^2 \quad (1)$$

where  $G_{i,j,k}(\mathbf{y}_i) = \frac{1}{2\pi\sigma^2} e^{-[(y_k^1 - i)^2 + (y_k^2 - j)^2]/2\sigma^2}$  is a Gaussian centred at joint  $y_k$  with fixed  $\sigma$ .

**Discussion.** As noted by [36], regressing coordinates directly is a highly non-linear and more difficult to learn mapping, which we also confirm here (Sect 5). The benefits of regressing a heatmap rather than  $(x, y)$  coordinates are twofold: first, one can understand failures and visualise the ‘thinking process’ of the network (see Figs 3 and 5); second, since by design, the output of the network can be multi-modal, *i.e.* allowed to have confidence at multiple spatial locations, learning becomes easier: early on in training (as shown in Fig 3), multiple locations may fire for a given joint; the incorrect ones are then slowly suppressed as training proceeds. In contrast, if the output were only the wrist  $(x, y)$  coordinate, the net would only have a lower loss if it gets its prediction right (even if it was ‘growing confidence’ in the correct position).

**Architecture.** The network architecture is shown in Fig 1, and sample activations for the layers are shown in Fig 5. To maximise the spatial resolution of the heatmap we make two important design choices: (i) minimal pooling is used (only two  $2 \times 2$  max-pooling layers), and (ii) all strides are

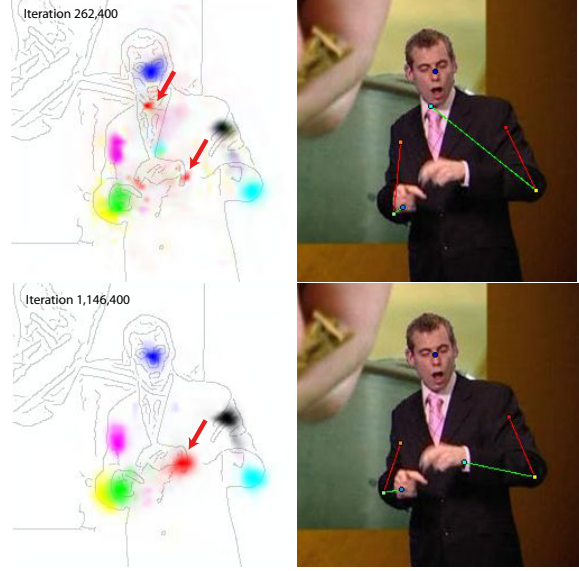


Figure 3. **Multiple peaks are possible with the Spatial ConvNet.** Early on in training (top), multiple locations may fire for a given joint. These are then suppressed as training proceeds (bottom). The arrows identify two modes for the wrist; one correct, one erroneous. As the training proceeds the erroneous one is diminished.

unity (so that the resolution is not reduced). All layers are followed by ReLUs except conv9 (the pooling layer). In contrast to AlexNet [21], our network is fully convolutional (no fully-connected layers) with the fully-connected layers of [21] replaced by  $1 \times 1$  convolutions. In contrast to both AlexNet and [35], our network is deeper, does not use local contrast normalisation (as we did not find this beneficial), and utilises less max-pooling.

## 2.2. Spatial fusion layers

Vanilla heatmap pose nets do not learn spatial dependencies of joints, and thus often predict kinematically impossible poses (see examples in the extended arXiv version of this paper). To address this, we add what we term ‘spatial fusion layers’ to the network. These spatial fusion layers (normal convolutional layers) take as an input pre-heatmap activations (conv7), and learn dependencies between the human body parts locations represented by these activations. In detail, these layers take as an input a concatenation of conv7 and conv3 (a skip layer), and feed these through five more convolutional layers with ReLUs (see Fig 4). Large kernels are used to inflate the receptive field of the network. We attach a separate loss layer to the end of this network and backpropagate through the whole network.

## 2.3. Optical flow for pose estimation

Given the heatmaps from the Spatial ConvNet from multiple frames, the heatmaps are reinforced with optical flow. This is done in three steps: (1) the confidences from nearby

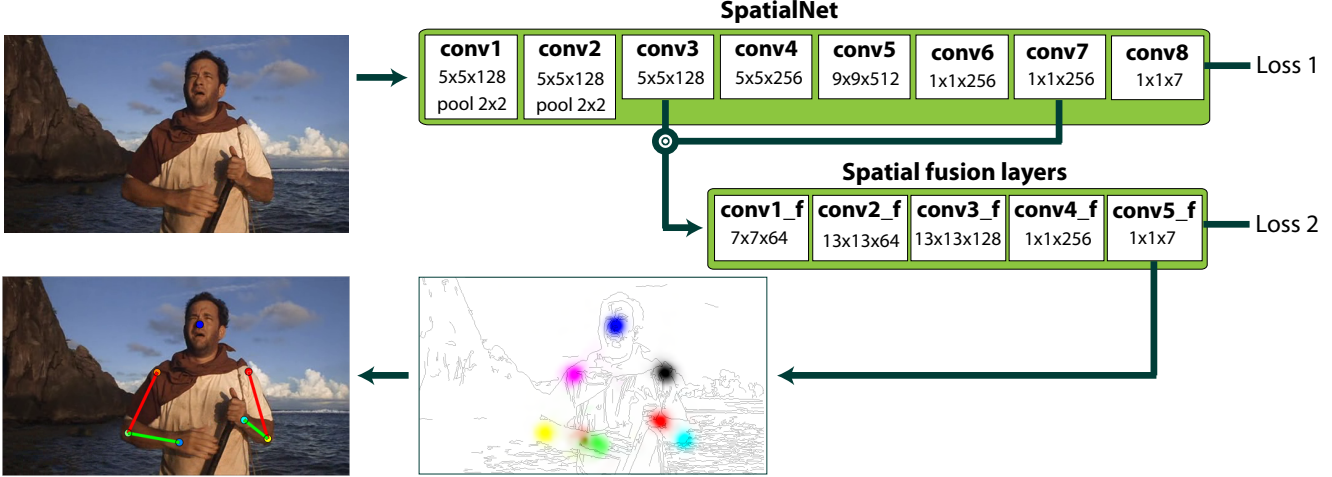


Figure 4. **Spatial fusion layers.** The fusion layers learn to encode dependencies between human body parts locations, learning an implicit spatial model.



Figure 5. **Sample activations for convolutional layers.** Neuron activations are shown for three randomly selected channels for each convolutional layer (resized here to the same size), with the input (pre-segmented for visualisation purposes) shown above. Low down in the net, neurons are activated at edges in the image (e.g. conv1 and conv2); higher up, they start responding more clearly to body parts (conv6 onwards). The outputs in conv8 are shown for the right elbow, left shoulder and left elbow.

frames are aligned to the current frame using dense optical flow; (2) these confidences are then *pooled* into a composite confidence map using an additional convolutional layer; and (3) the final upper body pose estimate for a frame is

then simply the positions of maximum confidence from the composite map. Below we discuss the first two steps.

#### Step 1: Warping confidence maps with optical flow.

For a given frame  $t$ , pixel-wise temporal tracks are computed from all neighbouring frames within  $n$  frames from  $((t - n) \text{ to } (t + n))$  to frame  $t$  using dense optical flow [38]. These optical flow tracks are used to warp confidence values in neighbouring confidence maps to align them to frame  $t$  by effectively shifting confidences along the tracks [4]. Example tracks and the warping of wrist confidence values are shown in Fig 6.

#### Step 2: Pooling the confidence maps.

The output of Step 1 is a set of confidence maps that are warped to frame  $t$ . From these ‘expert opinions’ about the joint positions, the task is first to select a confidence for each pixel for each joint, and then to select one position for each joint. One solution would be to simply average the warped confidence maps. However, not all experts should be treated equally: intuitively, frames further away (thus with more space for optical flow errors) should be given lower weight.

To this end we learn a parametric cross-channel pooling layer that takes as an input a set of warped heatmaps for a given joint, and as an output predicts a single ‘composite heatmap’. The input to this layer is a  $i \times j \times t$  heatmap volume, where  $t$  is the number of warped heatmaps (e.g. 31 for a neighbourhood of  $n = 15$ ). As the pooling layer, we train a  $1 \times 1$  kernel size convolutional layer for each joint. This is equivalent to cross-channel weighted sum-pooling, where we learn a single weight for each input channel (which correspond to the warped heatmaps). In total, we therefore learn  $t \times k$  weights (for  $k$  joints).



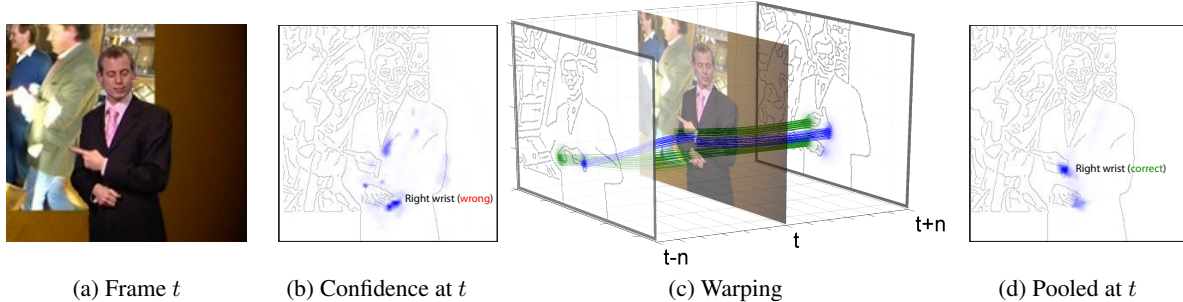


Figure 6. **Warping neighbouring heatmaps for improving pose estimates.** (a) RGB input at frame  $t$ . (b) Heatmap at frame  $t$  for right hand in the image. (c) Heatmaps from frames  $(t - n)$  and  $(t + n)$  warped to frame  $t$  using tracks from optical flow (green & blue lines). (d) Pooled confidence map with corrected modes.

	BBC	Ext. BBC	ChaLearn	PiW
Train frames	1.5M	7M	1M	4.5K (FLIC)
Test frames	1,000	1,000	3,200	830
Train labels	[2]	[2, 3]	Kinect	Manual
Test labels	Manual	Manual	Kinect	Manual
Train videos	13	85	393	-
Val videos	2	2	287	-
Test videos	5	5	275	30

Table 1. **Dataset overview, including train/val/test splits.**

### 3. Implementation Details

**Training.** The input frames are rescaled to height 256. A  $248 \times 248$  sub-image (of the  $N \times 256$  input image) is randomly cropped, randomly horizontally flipped, randomly rotated between  $-40^\circ$  and  $40^\circ$ , and resized to  $256 \times 256$ . Momentum is set to 0.95. The variance of the Gaussian is set to  $\sigma = 1.5$  with an output heatmap size of  $64 \times 64$ . A temporal neighbourhood of  $n = 15$  is input into the parametric pooling layer. The learning rate is set to  $10^{-4}$ , and decreased to  $10^{-5}$  at 80K iterations, to  $10^{-6}$  after 100K iterations and stopped at 120K iterations. We use Caffe [19].

**Training time.** Training was performed on four NVIDIA GTX Titan GPUs using a modified version of the Caffe framework [19] with multi-GPU support. Training SpatialNet on FLIC took 3 days & SpatialNet Fusion 6 days.

**Optical flow.** Optical flow is computed using FastDeepFlow [38] with Middlebury parameters.

### 4. Datasets

Experiments in this work are conducted on four large video pose estimation datasets, two from signed TV broadcasts, one of Italian gestures, and the third of Hollywood movies. An overview of the datasets is given in Table 1. The first three datasets are available at <http://www.robots.ox.ac.uk/~vgg/data/pose>.

**BBC Pose dataset.** This dataset [3] consists of 20 videos (each 0.5h–1.5h in length) recorded from the BBC with an

overlaid sign language interpreter. Each frame has been assigned pose estimates using the semi-automatic but reliable pose estimator of Buehler *et al.* [2] (used as training labels). 1,000 frames in the dataset have been manually annotated with upper-body pose (used as testing labels).

**Extended BBC Pose dataset.** This dataset [24] contains 72 additional training videos which, combined with the original BBC TV dataset, yields in total 85 training videos. The frames of these new videos have been assigned poses using the automatic tracker of Charles *et al.* [3]. The output of this tracker is noisier than the semi-automatic tracker of Buehler *et al.*, which results in partially noisy annotations.

**ChaLearn dataset.** The ChaLearn 2013 Multi-modal gesture dataset [9] contains 23 hours of Kinect data of 27 people. The data includes RGB, depth, foreground segmentations and full body skeletons. In this dataset, both the training and testing labels are noisy (from Kinect). The large variation in clothing across videos poses a challenging task for pose estimation methods.

**Poses in the Wild (PiW) and FLIC datasets.** The Poses in the Wild dataset [6] contains 30 sequences (total 830 frames) extracted from Hollywood movies. The frames are annotated with upper-body poses. It contains realistic poses in indoor and outdoor scenes, with background clutter, severe camera motion and occlusions. For training, we follow [6] and use all the images annotated with upper-body parts (about 4.5K) in the FLIC dataset [26].

### 5. Experiments

We first describe the evaluation protocol, then present comparisons to alternative network architectures, and finally give a comparison to state of the art. A demo video is online at <http://youtu.be/pj2N5DqBOgQ>.

#### 5.1. Evaluation protocol and details

**Evaluation protocol.** In all pose estimation experiments we compare the estimated joints against frames with manual

ground truth (except ChaLearn, where we compare against output from Kinect). We present results as graphs that plot accuracy vs distance from ground truth in pixels, where a joint is deemed correctly located if it is within a set distance of  $d$  pixels from a marked joint centre in ground truth.

**Experimental details.** All frames of the training videos are used for training (with each frame randomly augmented as detailed above). The frames are randomly shuffled prior to training to present maximally varying input data to the network. The hyperparameters (early stopping, variance  $\sigma$  etc.) are estimated using the validation set.

**Baseline method.** As a baseline method we include a CoordinateNet (described in [23]). This is a network with similar architecture to [28], but trained for regressing the joint positions directly (instead of a heatmap) [24].

**Computation time.** Our method is real-time (50fps on 1 GPU without optical flow, 5fps with optical flow).

## 5.2. Component evaluation

For these experiments the SpatialNet and baseline are trained and tested on the BBC Pose and Extended BBC Pose datasets. Fig 7 shows the results for wrists

With the SpatialNet, we observe a significant boost in performance (an additional 6.6%, from 79.6% to 86.1% at  $d = 6$ ) when training on the larger Extended BBC dataset compared to the BBC Pose dataset. As noted in Sect 4, this larger dataset is somewhat noisy. In contrast, the CoordinateNet is unable to make effective use of this additional noisy training data. We believe this is because its target (joint coordinates) does not allow for multi-modal output, which makes learning from noisy annotation challenging.

We observe a further boost in performance from using optical flow to warp heatmaps from neighbouring frames (an improvement of 2.6%, from 86.1% to 88.7% at  $d = 6$ ). Fig 9 shows the automatically learnt pooling weights. We see that for this dataset, as expected, the network learns to weigh frames temporally close to the current frame higher (because they contain less errors in optical flow).

Fig 8 shows a comparison of different pooling types (for cross-channel pooling). We compare learning a parametric pooling function to sum-pooling and to max-pooling (max-out [14]) across channels. As expected, parametric pooling performs best, and improves as the neighbourhood  $n$  increases. In contrast, results with both sum-pooling and max-pooling deteriorate as the neighbourhood size is increased further, as they are not able to down-weight predictions that are further away in time (and thus more prone to errors in optical flow). As expected, this effect is particularly noticeable for max-pooling.

**Failure modes.** The main failure mode for the vanilla heatmap network (conv1-conv8) occurs when multiple modes are predicted and the wrong one is selected (and the

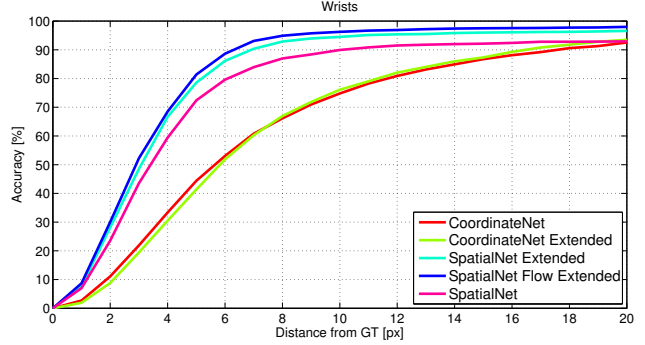


Figure 7. **Comparison of the performance of our nets for wrists on BBC Pose.** Plots show accuracy as the allowed distance from manual ground truth is increased. CoordinateNet is the network in [23]; SpatialNet is the heatmap network; and SpatialNet Flow is the heatmap network with the parametric pooling layer. ‘Extended’ indicates that the network is trained on Extended BBC Pose instead of BBC Pose. We observe a significant gain for the SpatialNet from using the additional training data in the Extended BBC dataset (automatically labelled – see Sect 4) training data, and a further boost from using optical flow information (and selecting the warping weights with the parametric pooling layer).

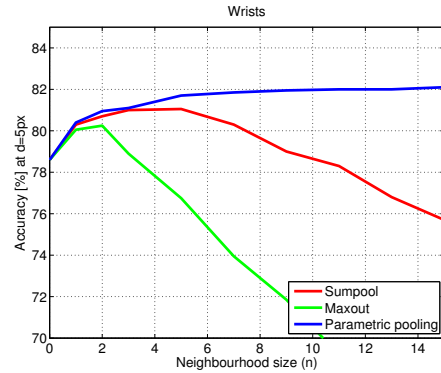


Figure 8. **Comparison of pooling types.** Results are shown for wrists in BBC Pose at threshold  $d = 5$ px. Parametric pooling (learnt cross-channel pooling weights) performs best.

resulting poses are often kinematically impossible for a human to perform). Examples of these failures are included in the extended arXiv version of this paper. The spatial fusion layers resolve these failures.

## 5.3. Comparison to state of the art

**Training.** We investigated a number of strategies for training on these datasets including training from scratch (using only the training data provided with the dataset), or training on one (*i.e.* BBC Pose) and fine-tuning on the others. We found that provided the first and last layers of the Spatial Net are initialized from (any) trained heatmap network, the rest can be trained either from scratch or fine-tuned with similar performance. We hypothesise this is because the datasets are very different – BBC Pose contains long-sleeved persons, ChaLearn short-sleeved persons and

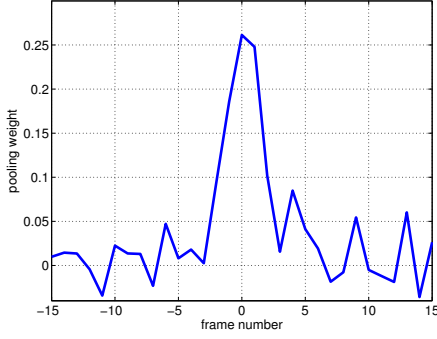


Figure 9. **Learnt pooling weights for BBC Pose with  $n = 15$ .** Weights shown for the right wrist. The centre frame receives highest weight. The jitter in weights is due to errors in optical flow computation (caused by the moving background in the video) – the errors become larger further away from the central frame (hence low or even negative weights far away).

Poses in the Wild contains non-frontal poses with unusual viewing angles. For all the results reported here we train BBC Pose from scratch, initialize the first and last layer from this, and fine-tune on training data of other datasets.

**BBC Pose.** Fig 10 shows a comparison to the state of the art on the BBC Pose dataset. We compare against all previous reported results on the dataset. These include Buehler *et al.* [2], whose pose estimator is based on a pictorial structure model; Charles *et al.* (2013) [3] who uses a Random Forest; Charles *et al.* (2014) [4] who predict joints sequentially with a Random Forest; Pfister *et al.* (2014) [24] who use a deep network similar to our CoordinateNet (with multiple input frames); and the deformable part-based model of Yang & Ramanan (2013) [39].

We outperform all previous work by a large margin, with a particularly noticeable gap for wrists (an addition of 10% compared to the best competing method at  $d = 6$ ).

**Chalearn.** Fig 11 shows a comparison to the state of the art on ChaLearn. We again outperform the state of the art even without optical flow (an improvement of 3.5% at  $d = 6$ ), and observe a further boost by using optical flow (beating state of the art by an addition of 5.5% at  $d = 6$ ), and a significant further improvement from using a deeper network (an additional 13% at  $d = 6$ ).

**Poses in the Wild.** Figs 12 & 14 show a comparison to the state of the art on Poses in the Wild. We replicate the results of the previous state of the art method using code provided by the authors [6]. We outperform the state of the art on this dataset by a large margin (an addition of 30% for wrists and 24% for elbows at  $d = 8$ ). Using optical flow yields a significant 10% improvement for wrists and 13% for elbows at  $d = 8$ . Fig 15 shows example predictions.

**FLIC.** Fig 13 shows a comparison to the state of the art on FLIC. We outperform all pose estimation methods that

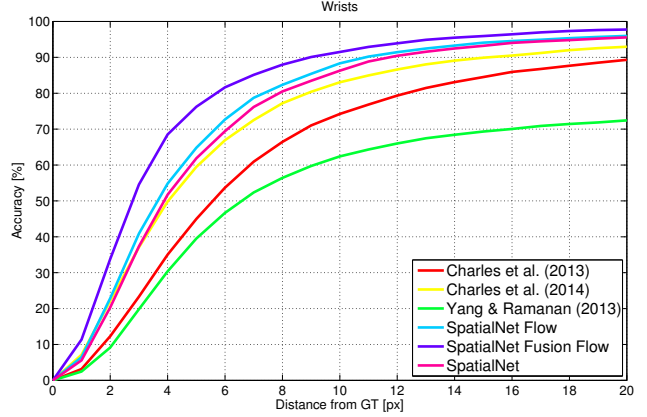


Figure 11. **Comparison to the state of the art on ChaLearn.** Our method outperforms state of the art by a large margin (an addition of 19% at  $d = 4$ ). See arXiv version for elbows & shoulders.

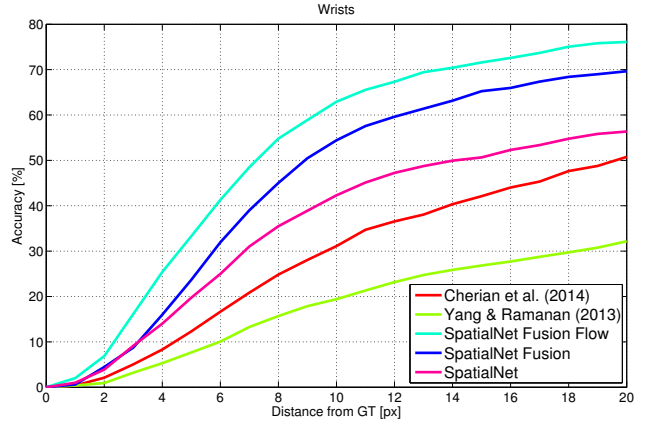


Figure 12. **Comparison to state of the art on Poses in the Wild.** Our method outperforms state of the art by a large margin (an addition of 30% at  $d = 8$ , with 10% from flow).

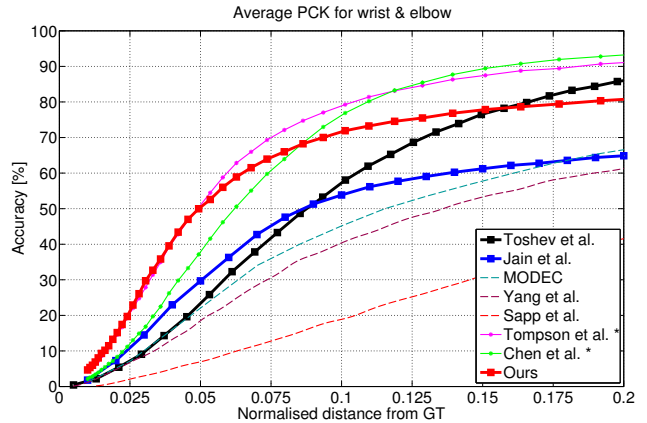


Figure 13. **Comparison to state of the art on FLIC.** Solid lines represent deep models; methods with a square (■) are without a graphical model; methods with an asterisk (\*) are with a graphical model. Our method outperforms competing methods without a graphical model by a large margin in the high precision area (an addition of 20% at  $d = 0.05$ ).

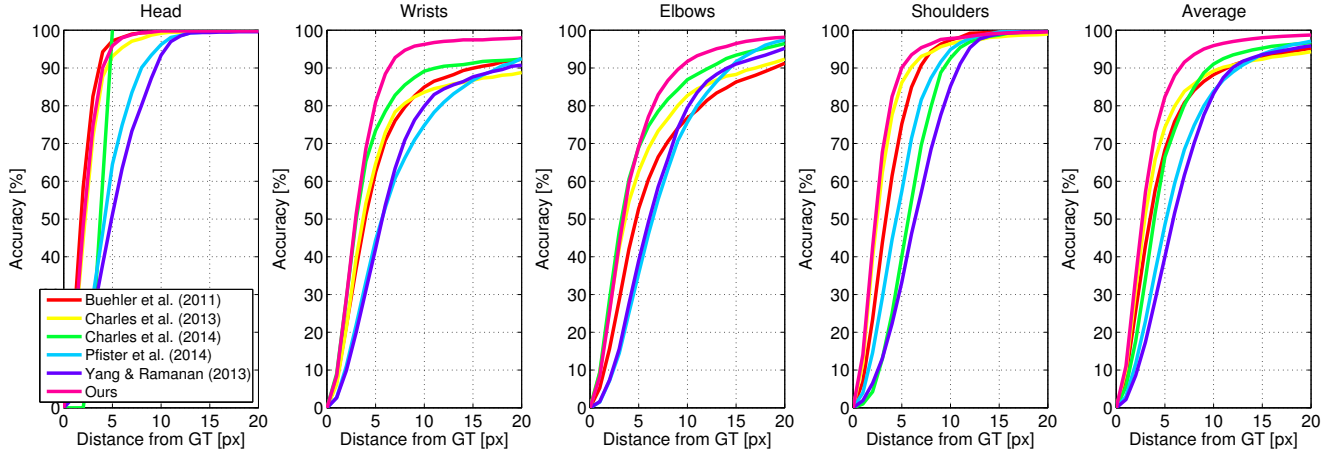


Figure 10. **Comparison to the state of the art on BBC Pose.** Plots show accuracy per joint type (average over left and right body parts) as the allowed distance from manual ground truth is increased. We outperform all previous work by a large margin; notice particularly the performance for wrists, where we outperform the best competing method with an addition of 10% at  $d = 6$ . Our method uses SpatialNet Flow Extended. Pfister *et al.* (2014) uses Extended BBC Pose; Buehler *et al.*, Charles *et al.* and Yang & Ramanan use BBC Pose.

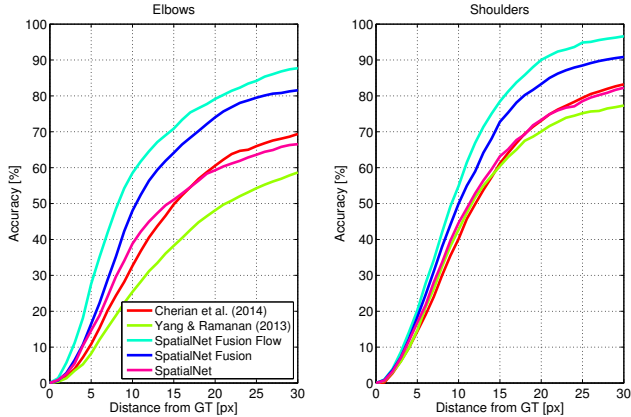


Figure 14. **Poses in the Wild: elbows & shoulders.**

don't use a graphical model, and match or even slightly outperform graphical model-based methods [5, 35] in the very high precision region ( $< 0.05$  from GT). The increase in accuracy at  $d = 0.05$  is 20% compared to methods not using a graphical model, and 12% compared to [5] who use a graphical model. Thompson *et al.* is [35]; Jain *et al.* is [17]. Predictions are provided by the authors of [5, 35] and evaluation code by the authors of [35].

## 6. Conclusion

We have presented a new architecture for pose estimation in videos that is able to utilize appearances across multiple frames. The proposed ConvNet is a simple, direct method for regressing heatmaps, and its performance is improved by combining it with optical flow and spatial fusion layers. We have also shown that our method outperforms the state of the art on three large video pose estimation datasets.

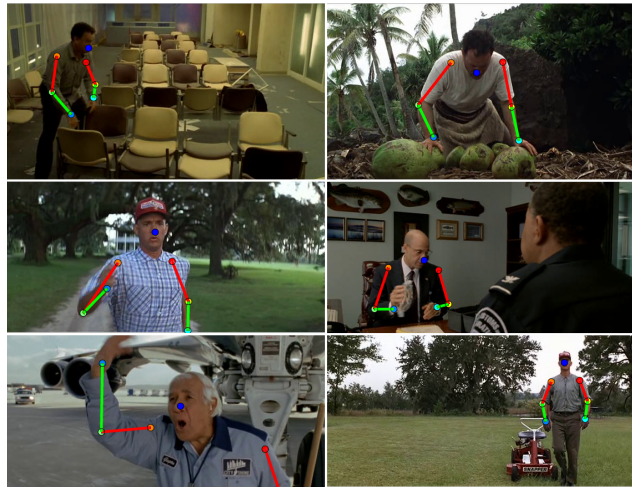


Figure 15. **Example predictions on Poses in the Wild.**

Further improvements may be obtained by using additional inputs for the spatial ConvNet, for example multiple RGB frames [24] or optical flow [18] – although prior work has shown little benefit from this so far.

The benefits of aligning pose estimates from multiple frames using optical flow, as presented here, are complementary to architectures that explicitly add spatial MRF and refinement layers [35, 36].

Finally, we have demonstrated the architecture for human pose estimation, but a similar optical flow-mediated combination of information could be used for other tasks in video, including classification and segmentation.

**Acknowledgements:** Financial support was provided by Osk. Huttunen Foundation and EPSRC grant EP/I012001/1.



## References

- [1] L. Bourdev and J. Malik. Poselets: Body part detectors trained using 3d human pose annotations. In *Proc. CVPR*, 2009. [1](#)
- [2] P. Buehler, M. Everingham, D. P. Huttenlocher, and A. Zisserman. Upper body detection and tracking in extended signing sequences. *IJCV*, 2011. [1](#), [5](#), [7](#)
- [3] J. Charles, T. Pfister, M. Everingham, and A. Zisserman. Automatic and efficient human pose estimation for sign language videos. *IJCV*, 2013. [1](#), [5](#), [7](#)
- [4] J. Charles, T. Pfister, D. Magee, D. Hogg, and A. Zisserman. Upper body pose estimation with temporal sequential forests. *Proc. BMVC*, 2014. [1](#), [2](#), [4](#), [7](#)
- [5] X. Chen and A. Yuille. Articulated pose estimation with image-dependent preference on pairwise relations. In *Proc. NIPS*, 2014. [1](#), [8](#)
- [6] A. Cherian, J. Mairal, K. Alahari, and C. Schmid. Mixing body-part sequences for human pose estimation. In *Proc. CVPR*, 2014. [5](#), [7](#)
- [7] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. *Proc. ICML*, 2014. [1](#)
- [8] M. Eichner, M. Marin-Jimenez, A. Zisserman, and V. Ferrari. 2d articulated human pose estimation and retrieval in (almost) unconstrained still images. *IJCV*, 2012. [1](#)
- [9] S. Escalera, J. Gonzalez, X. Baro, M. Reyes, O. Lopes, I. Guyon, V. Athistos, and H. Escalante. Multi-modal gesture recognition challenge 2013: Dataset and results. In *Proc. ICMI*, 2013. [5](#)
- [10] P. Felzenszwalb and D. Huttenlocher. Pictorial structures for object recognition. *IJCV*, 2005. [1](#)
- [11] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proc. CVPR*, 2014. [1](#)
- [12] R. Girshick, J. Shotton, P. Kohli, A. Criminisi, and A. Fitzgibbon. Efficient regression of general-activity human poses from depth images. In *Proc. ICCV*, 2011. [1](#)
- [13] G. Gkioxari, B. Hariharan, R. Girshick, and J. Malik. Using k-poselets for detecting people and localizing their keypoints. In *Proc. CVPR*, 2014. [1](#)
- [14] I. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks. *J. Machine Learning Research*, 2013. [6](#)
- [15] I. J. Goodfellow, Y. Bulatov, J. Ibarz, S. Arnoud, and V. Shet. Multi-digit number recognition from street view imagery using deep convolutional neural networks. In *Proc. ICLR*, 2014. [1](#)
- [16] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. *Proc. NIPS Workshops*, 2014. [1](#)
- [17] A. Jain, J. Tompson, M. Andriluka, G. Taylor, and C. Bregler. Learning human pose estimation features with convolutional networks. *Proc. ICLR*, 2014. [1](#), [8](#)
- [18] A. Jain, J. Tompson, Y. LeCun, and C. Bregler. MoDeep: A deep learning framework using motion features for human pose estimation. *Proc. ACCV*, 2014. [2](#), [8](#)
- [19] Y. Jia. Caffe: An open source convolutional architecture for fast feature embedding. <http://caffe.berkeleyvision.org/>, 2013. [5](#)
- [20] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proc. CVPR*, 2014. [1](#)
- [21] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. NIPS*, 2012. [1](#), [3](#)
- [22] M. Osadchy, Y. LeCun, and M. Miller. Synergistic face detection and pose estimation with energy-based models. *JMLR*, 2007. [1](#)
- [23] T. Pfister. *Advancing Human Pose and Gesture Recognition*. PhD thesis, University of Oxford, 2015. [6](#)
- [24] T. Pfister, K. Simonyan, J. Charles, and A. Zisserman. Deep convolutional neural networks for efficient pose estimation in gesture videos. *Proc. ACCV*, 2014. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#)
- [25] S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. CNN features off-the-shelf: an astounding baseline for recognition. *CVPR Workshops*, 2014. [1](#)
- [26] B. Sapp and B. Taskar. Modec: Multimodal decomposable models for human pose estimation. In *Proc. CVPR*, 2013. [5](#)
- [27] B. Sapp, D. Weiss, and B. Taskar. Parsing human motion with stretchable models. In *Proc. CVPR*, 2011. [1](#)
- [28] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *Proc. ICLR*, 2014. [1](#), [6](#)
- [29] J. Shotton, R. Girshick, A. Fitzgibbon, T. Sharp, M. Cook, M. Finocchio, R. Moore, P. Kohli, A. Criminisi, A. Kipman, and A. Blake. Efficient human pose estimation from single depth images. *PAMI*, 2013. [1](#)
- [30] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. *Proc. NIPS*, 2014. [1](#), [2](#)
- [31] M. Sun, P. Kohli, and J. Shotton. Conditional regression forests for human pose estimation. In *Proc. CVPR*, 2012. [1](#)
- [32] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proc. CVPR*, 2014. [1](#)
- [33] G. Taylor, R. Fergus, G. Williams, I. Spiro, and C. Bregler. Pose-sensitive embedding by nonlinear nca regression. In *Proc. NIPS*, 2010. [1](#)
- [34] J. Taylor, J. Shotton, T. Sharp, and A. Fitzgibbon. The vitruvian manifold: Inferring dense correspondences for one-shot human pose estimation. In *Proc. CVPR*, 2012. [1](#)
- [35] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler. Efficient object localization using convolutional networks. *Proc. CVPR*, 2015. [1](#), [3](#), [8](#)
- [36] J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. *Proc. NIPS*, 2014. [1](#), [3](#), [8](#)
- [37] A. Toshev and C. Szegedy. DeepPose: Human pose estimation via deep neural networks. *CVPR*, 2014. [1](#), [2](#)
- [38] P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid. Deepflow: Large displacement optical flow with deep matching. In *Proc. ICCV*, 2013. [4](#), [5](#)
- [39] Y. Yang and D. Ramanan. Articulated human detection with flexible mixtures of parts. *PAMI*, 2013. [1](#), [7](#)
- [40] M. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. *Proc. ECCV*, 2014. [1](#)