

Deep Learning for various CV tasks

Computer Vision (SJK02)

Universitat Jaume I

Part A:

Classification

Segmentation

Object detection

(Image-based) biometrics

Part B:

Sequence processing

Optical flow

Action Recognition

Self-supervised learning

Transformers

Sequence processing: RNNs, LSTMs

Recurrent Neural Networks

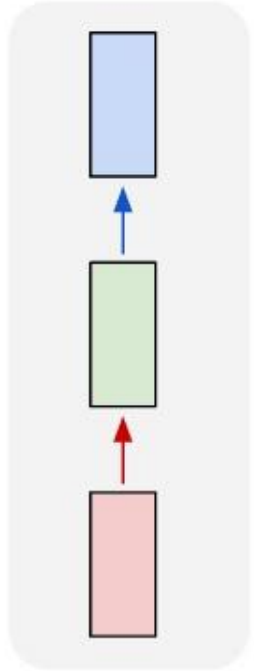


Unrolling

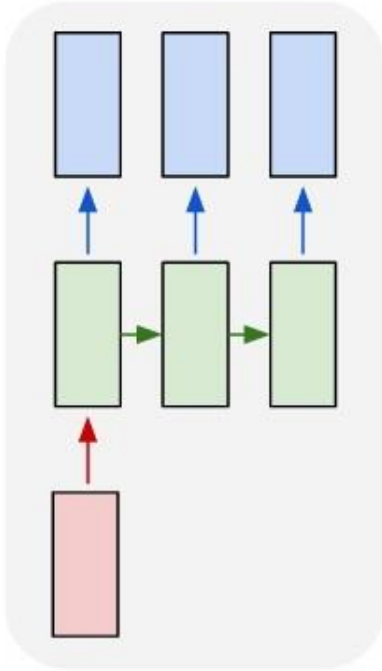


Flexible input/output sizes

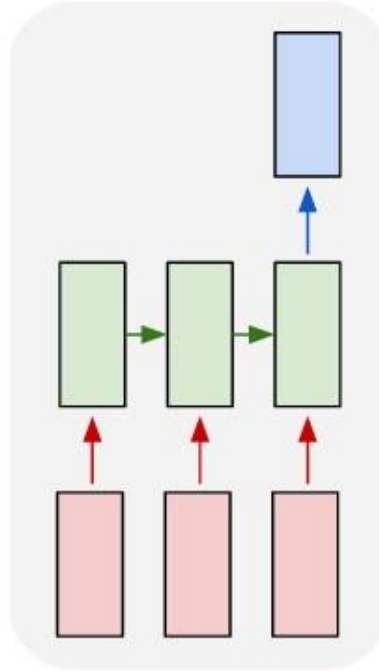
one to one



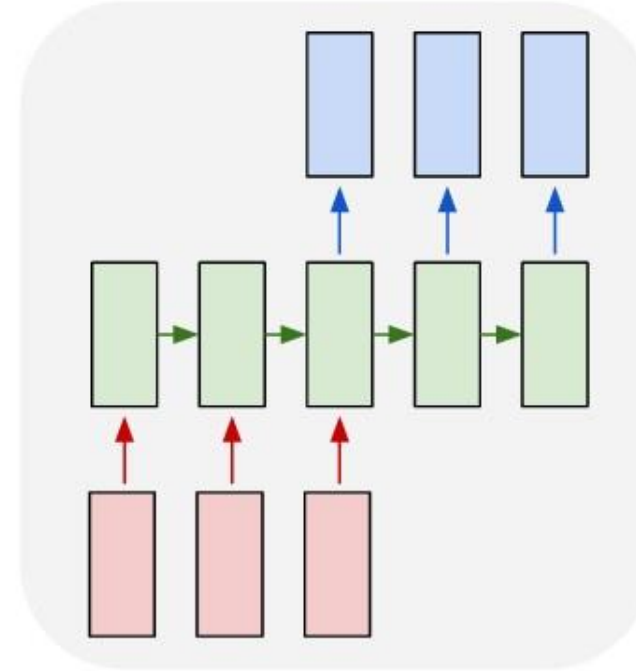
one to many



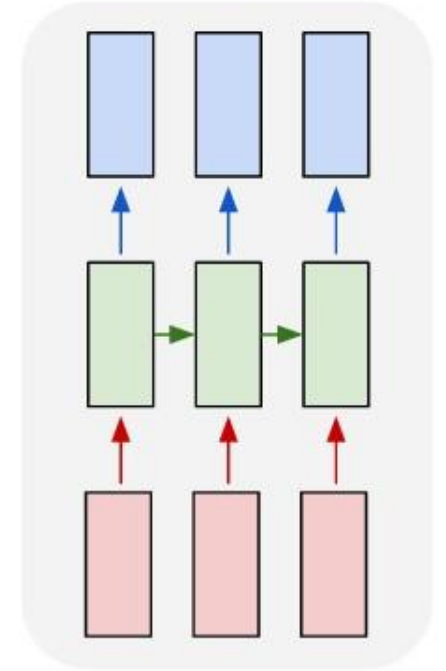
many to one



many to many



many to many



Match problems
to models

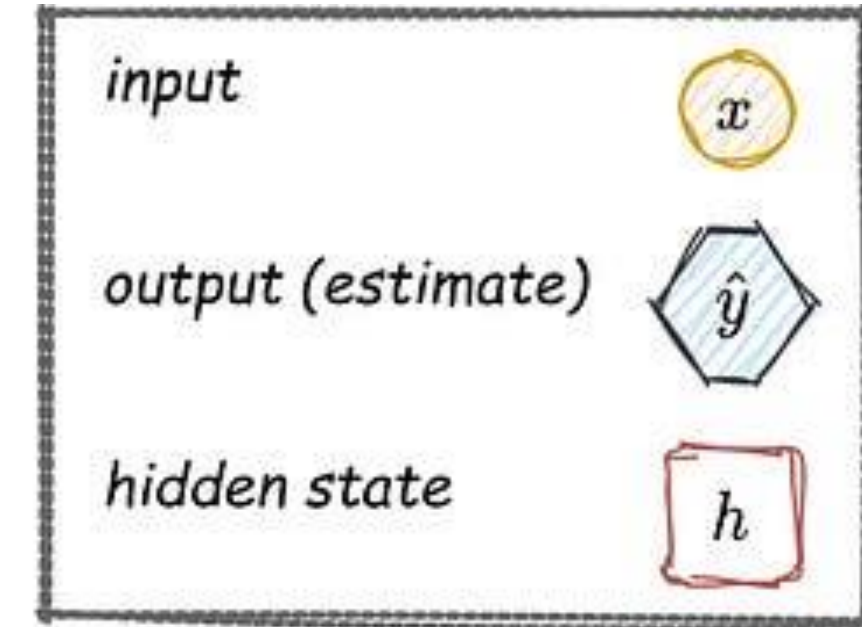
Frame-wise **video**
classification

Image captioning

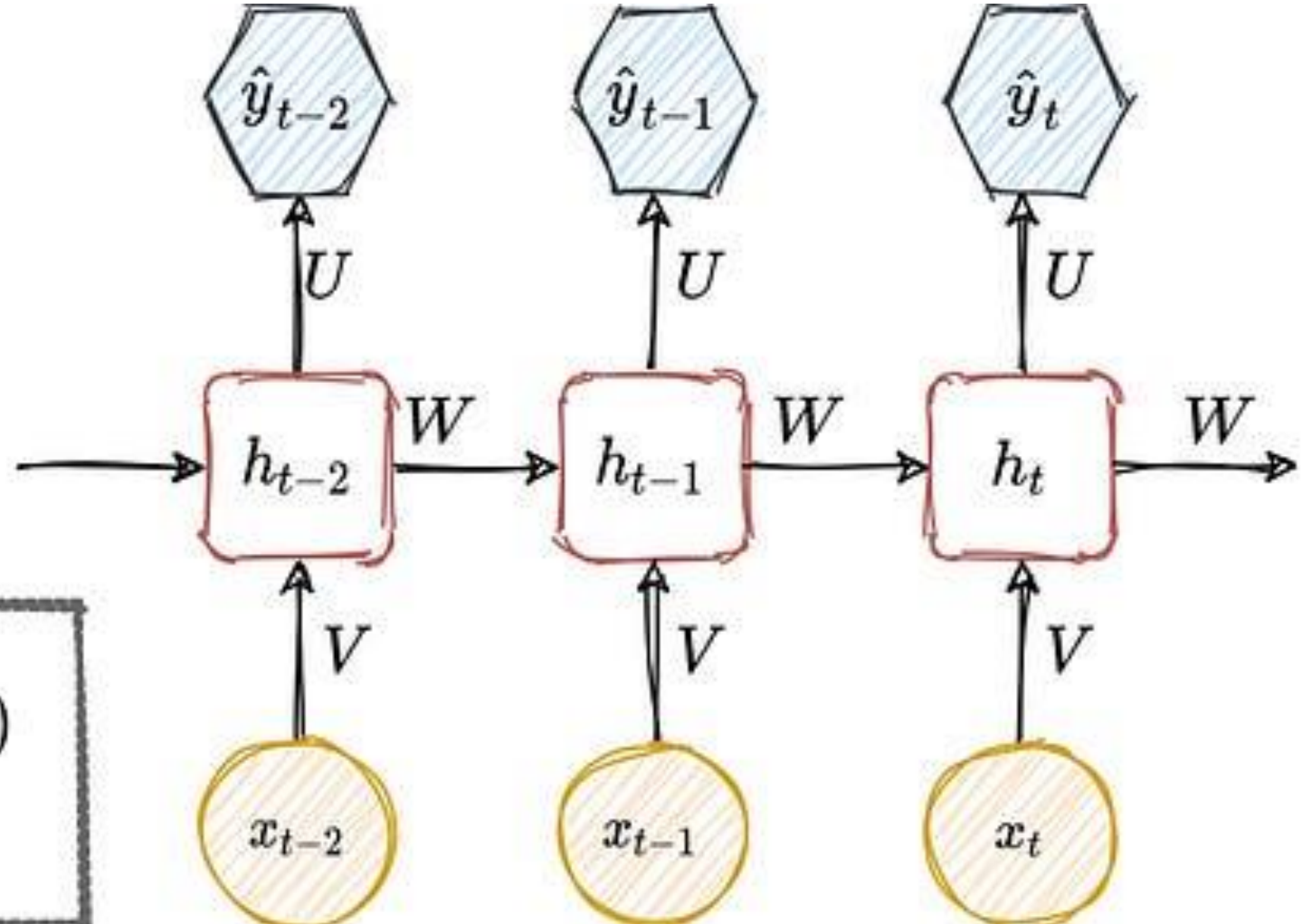
Image
classification

Sentiment analysis

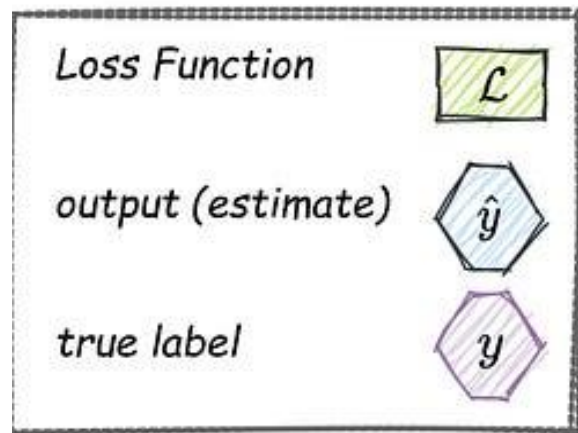
Machine translation



$$h_t = f_h (Vx_t + Wh_{t-1} + b_h)$$
$$\hat{y}_t = f_y (Uh_t + b_y)$$



Backpropagation through time (BTT)



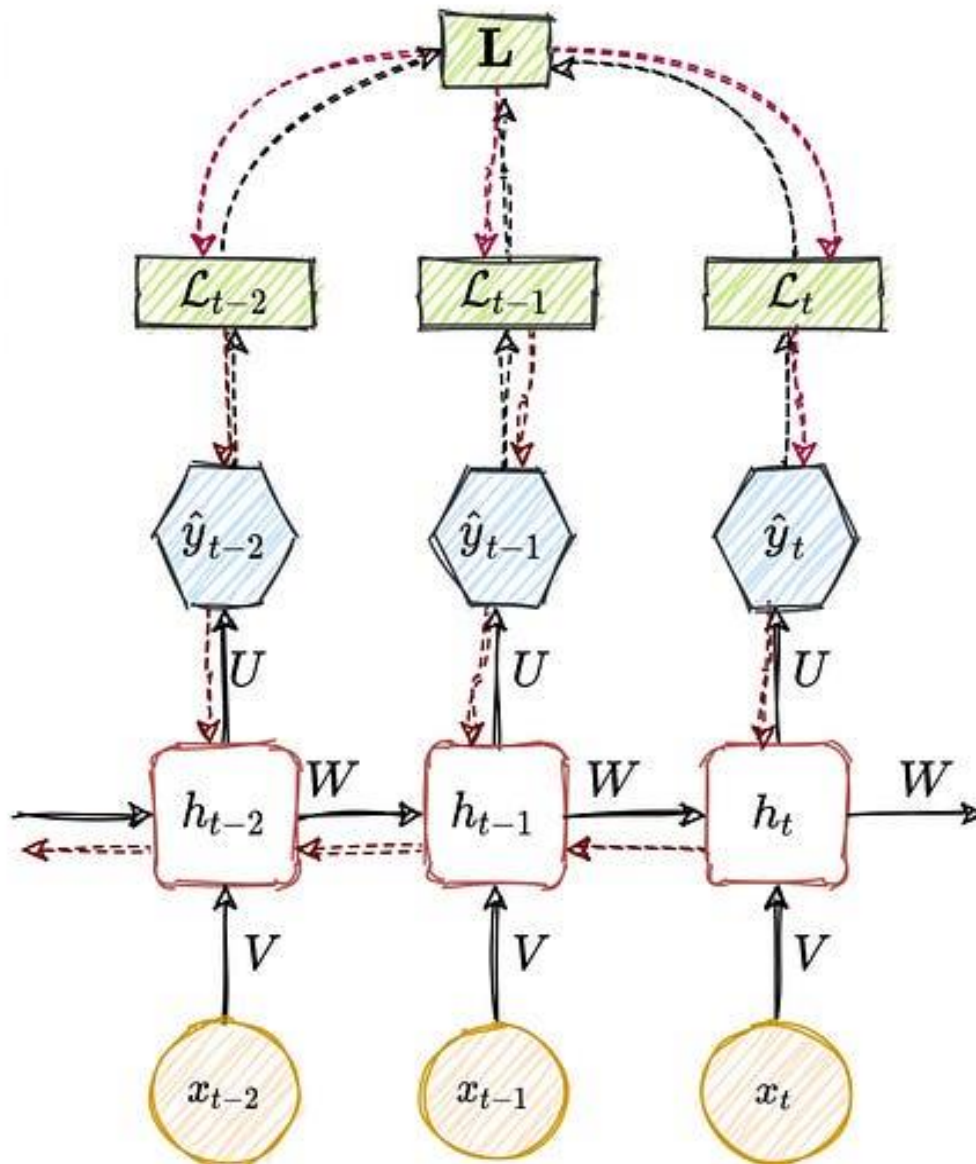
$$\mathbf{L} = \sum_i \mathcal{L}_i(\hat{y}_t, y_t)$$

Forward Pass:

$$h_t, \hat{y}_t, \mathcal{L}_t, \mathbf{L}$$

Backward Pass:

$$\frac{\partial \mathbf{L}}{\partial U}, \frac{\partial \mathbf{L}}{\partial V}, \frac{\partial \mathbf{L}}{\partial W}, \frac{\partial \mathbf{L}}{\partial b_h}, \frac{\partial \mathbf{L}}{\partial b_y}$$



Exploding/vanishing gradients

$$\frac{\partial \mathbf{L}}{\partial W} \propto \sum_{i=0}^T \left(\prod_{i=k+1}^y \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W}$$

1. *Vanishing gradient*

$$\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\|_2 < 1$$

2. *Exploding gradient*

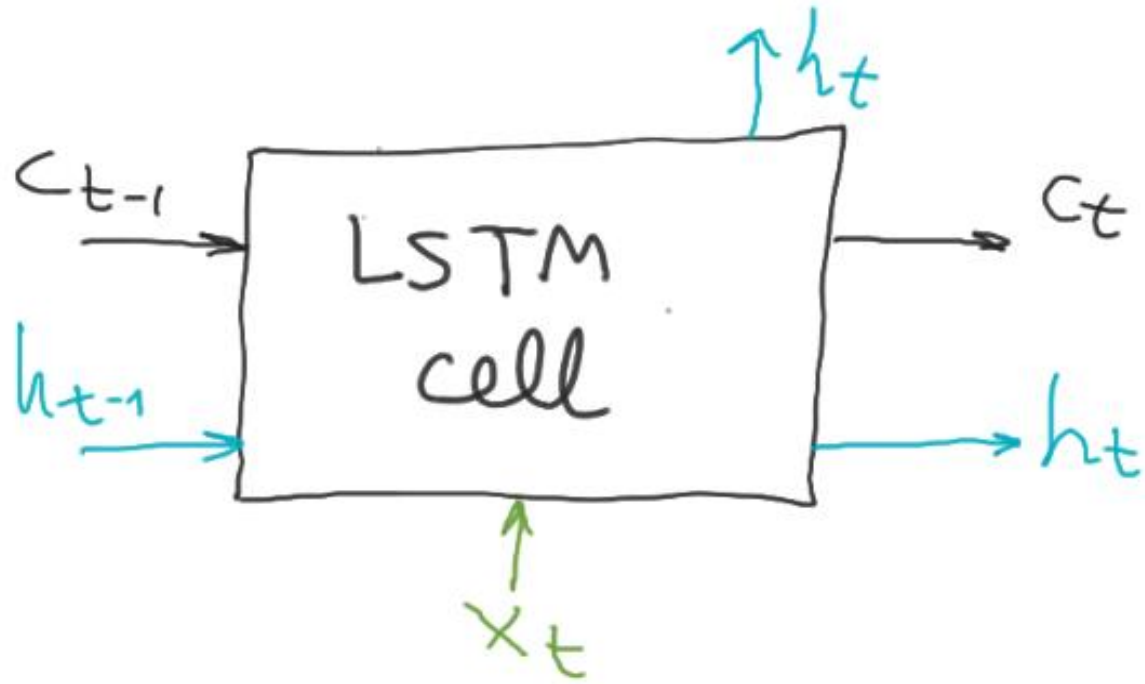
$$\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\|_2 > 1$$



Dealing with this problem

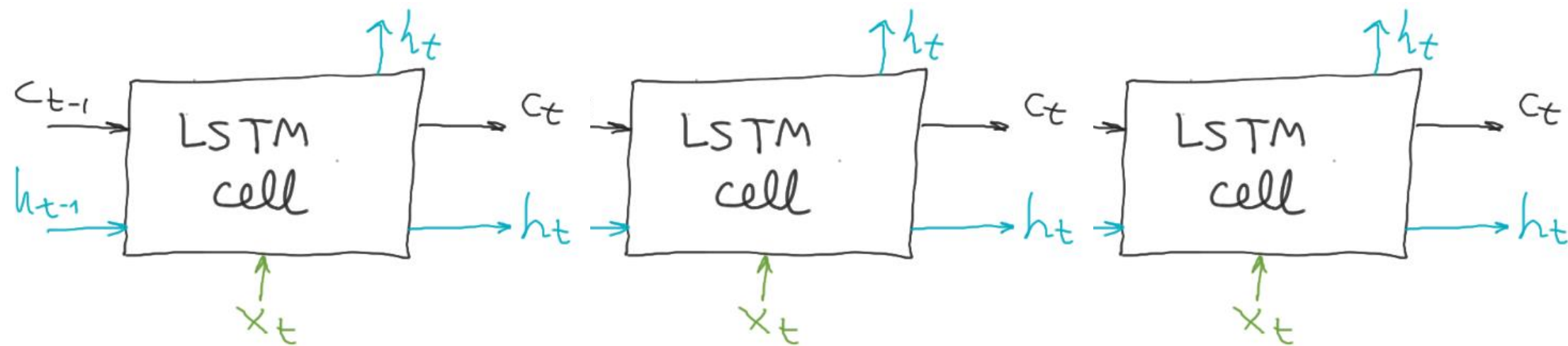
- Skip connections (like ResNets)
- Remove 1-length connections
- Leaky recurrent units (regulate how much passed by α)
- Gated RN (learnable α)
- **LSTMs**

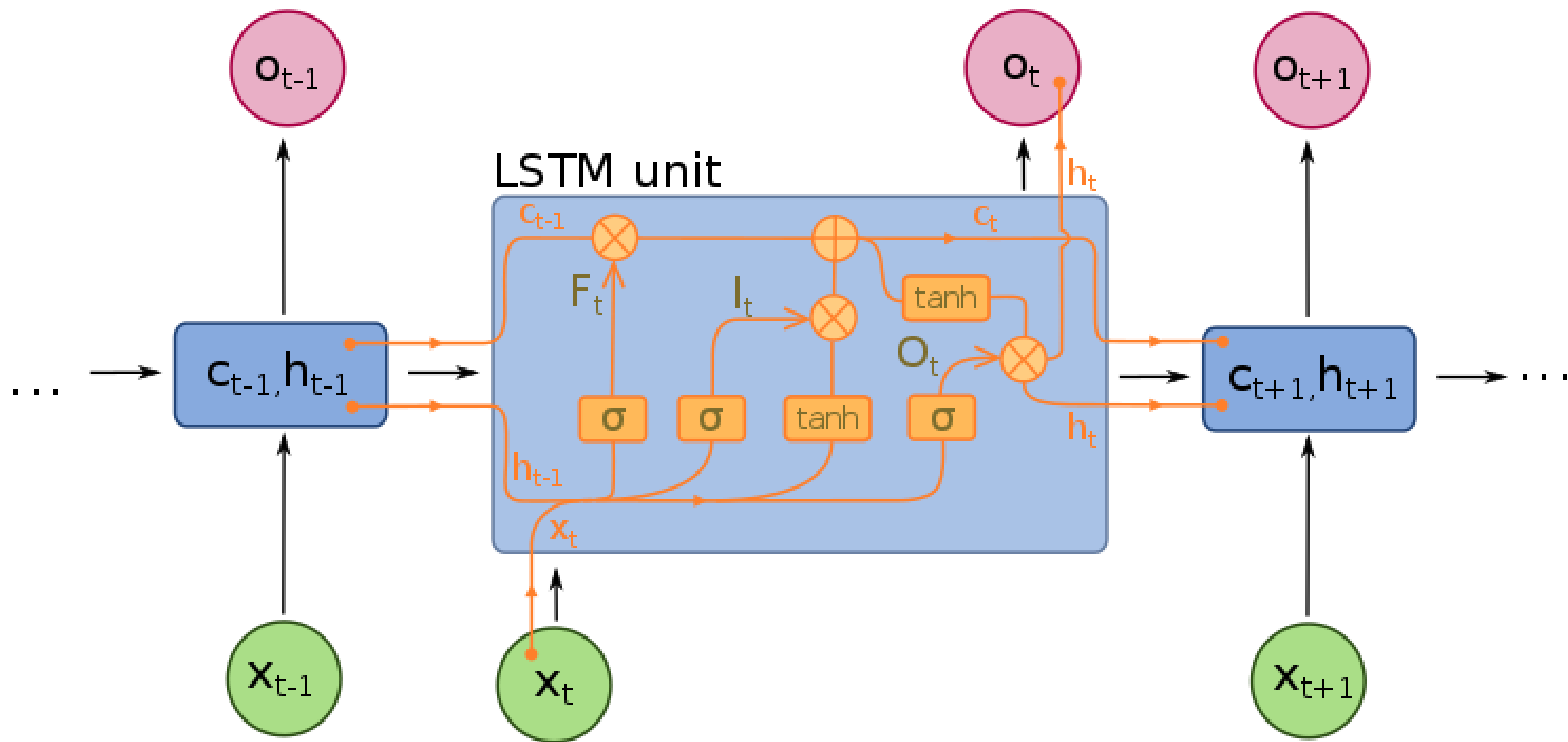
Long short-term memory (LSTM)

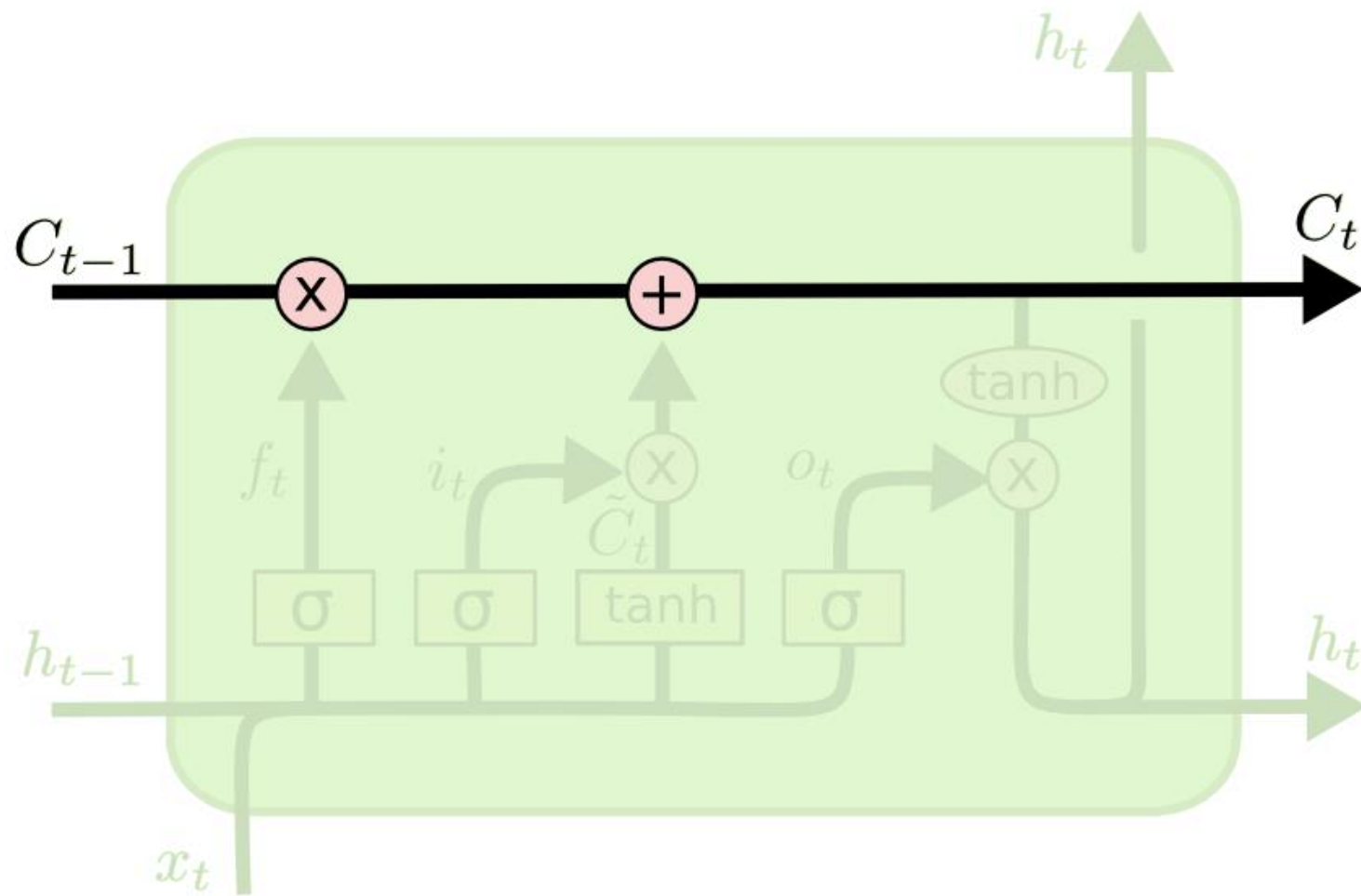


$$\left\{ \begin{array}{l} x_t = \text{input} \\ c_t = \text{cell state ("memory")} \\ h_t = \text{hidden state (= output)} \end{array} \right.$$

Unrolling

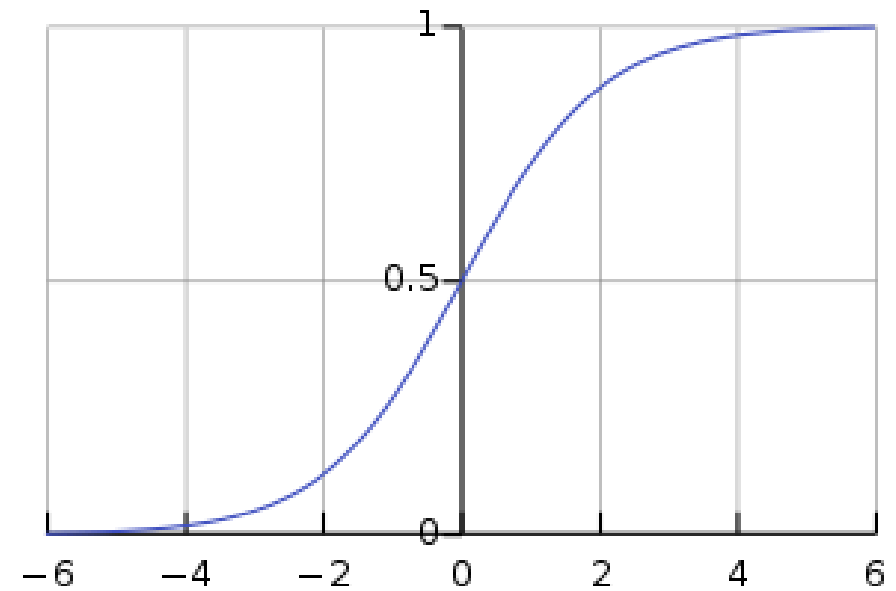
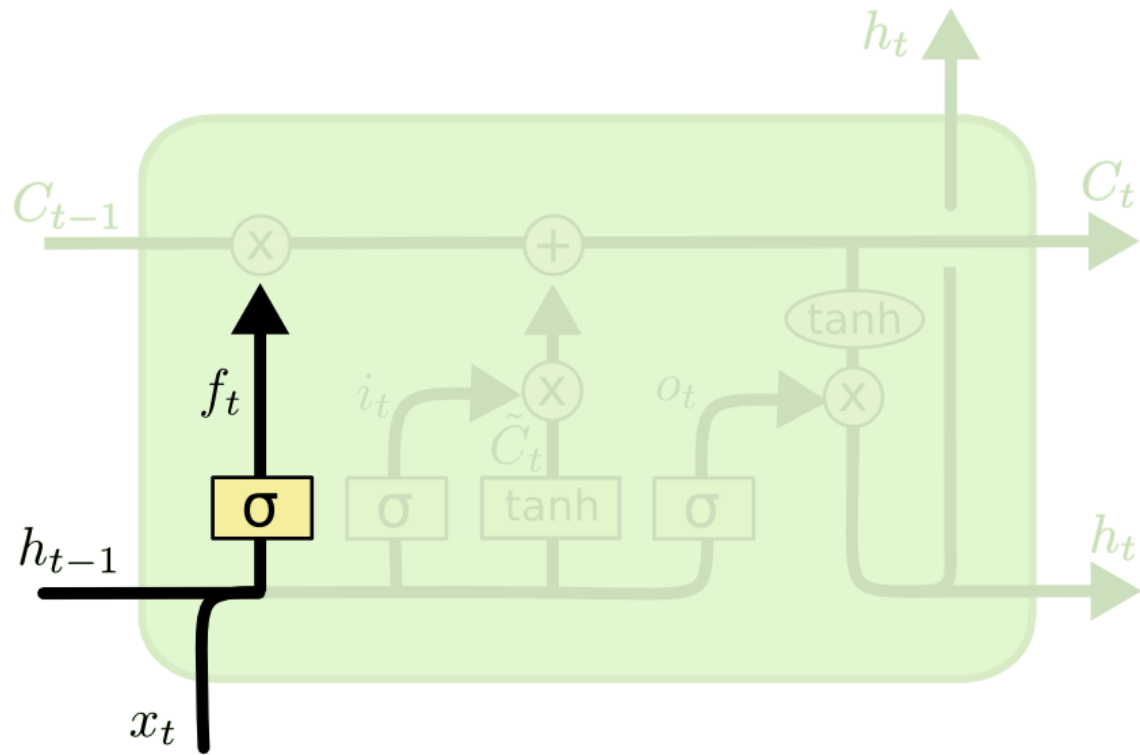






Cell **state** can go unchanged
Gates regulate how much it is changed

Forget gate

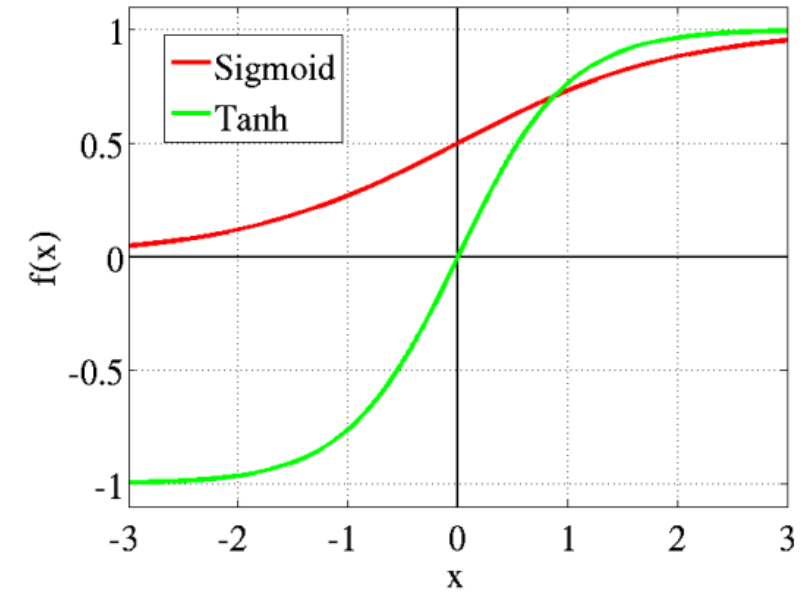
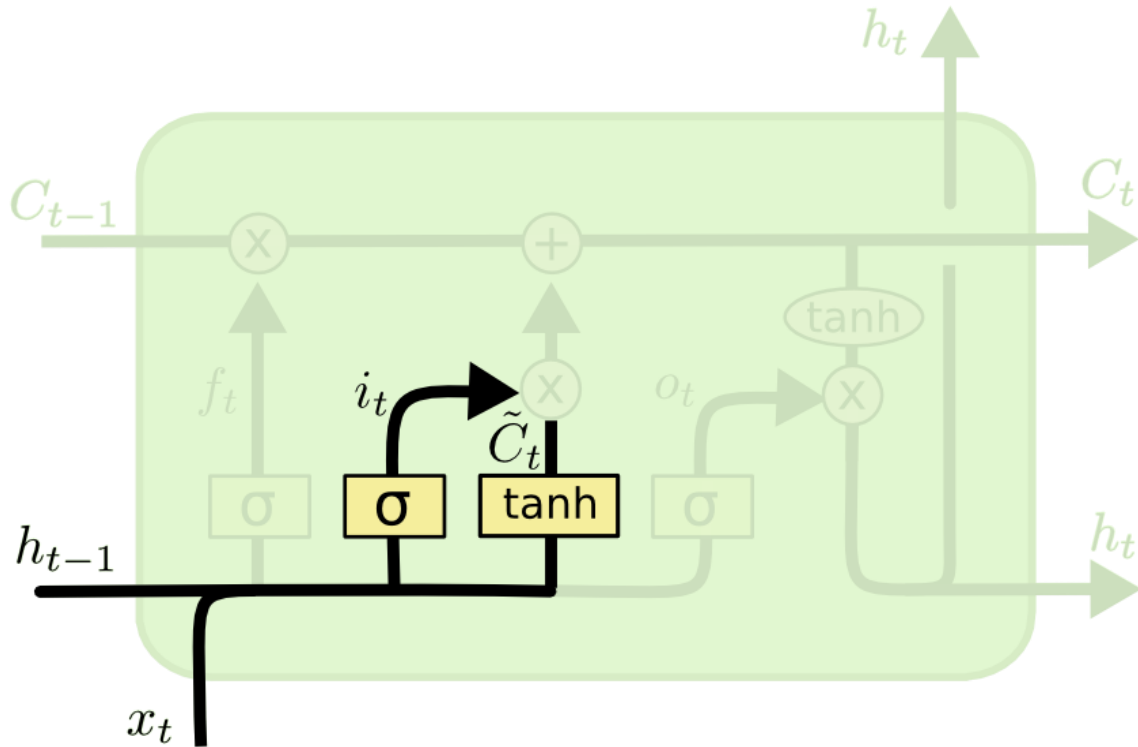


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x} = 1 - \sigma(-x).$$

How much, between 0 (nothing) and 1 (all), each element in C_{t-1} is kept

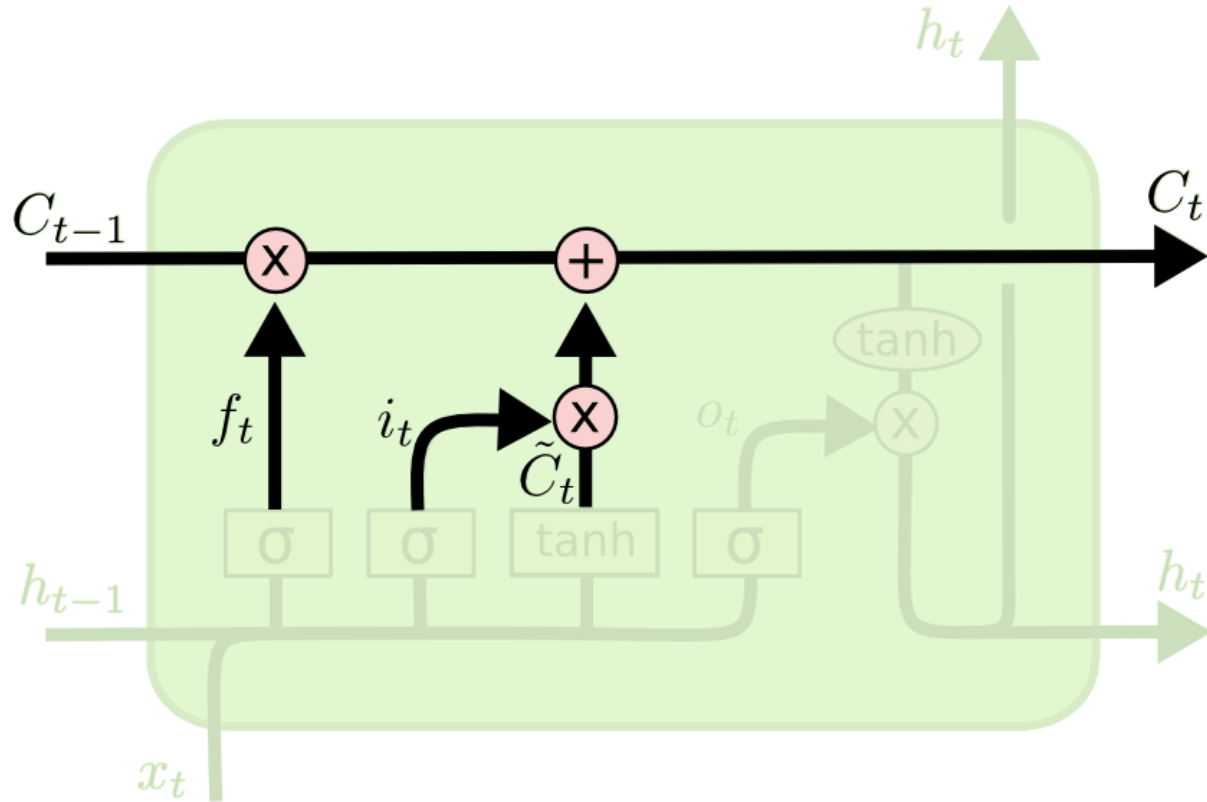
Input gate



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

New information and how much of it to add to C_{t-1}

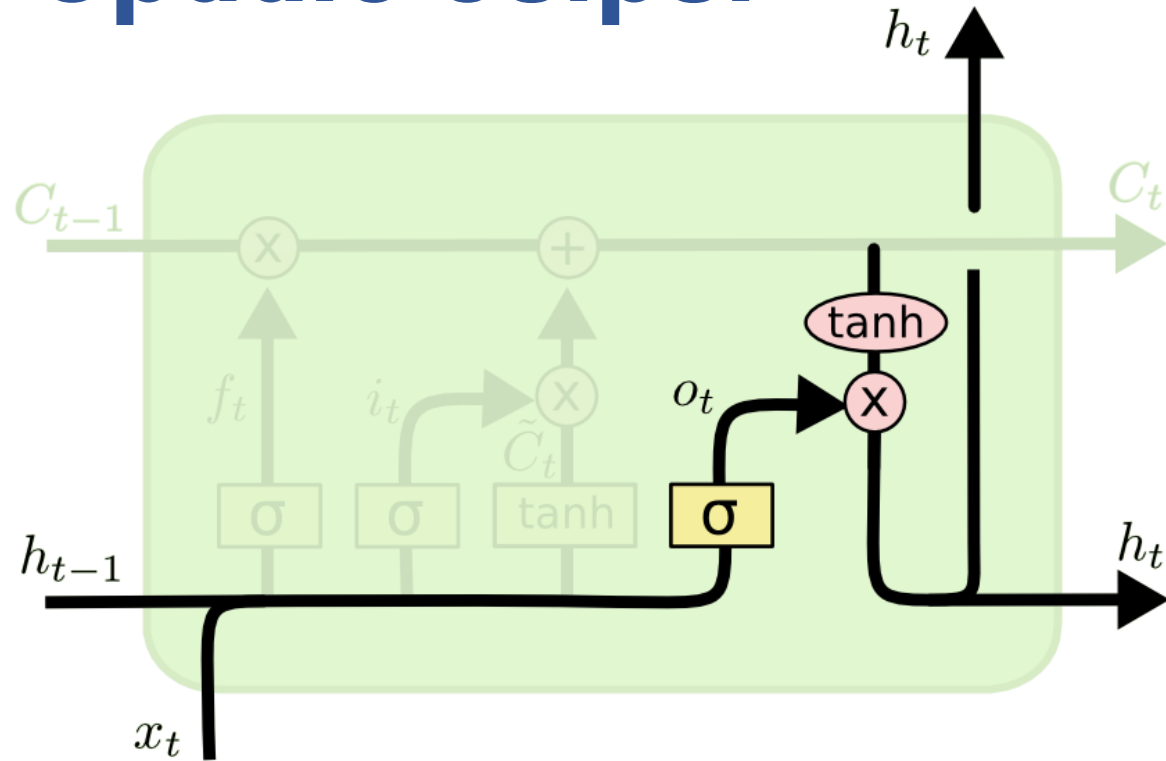
Update state



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Update C_{t-1} to get C_t

Update output

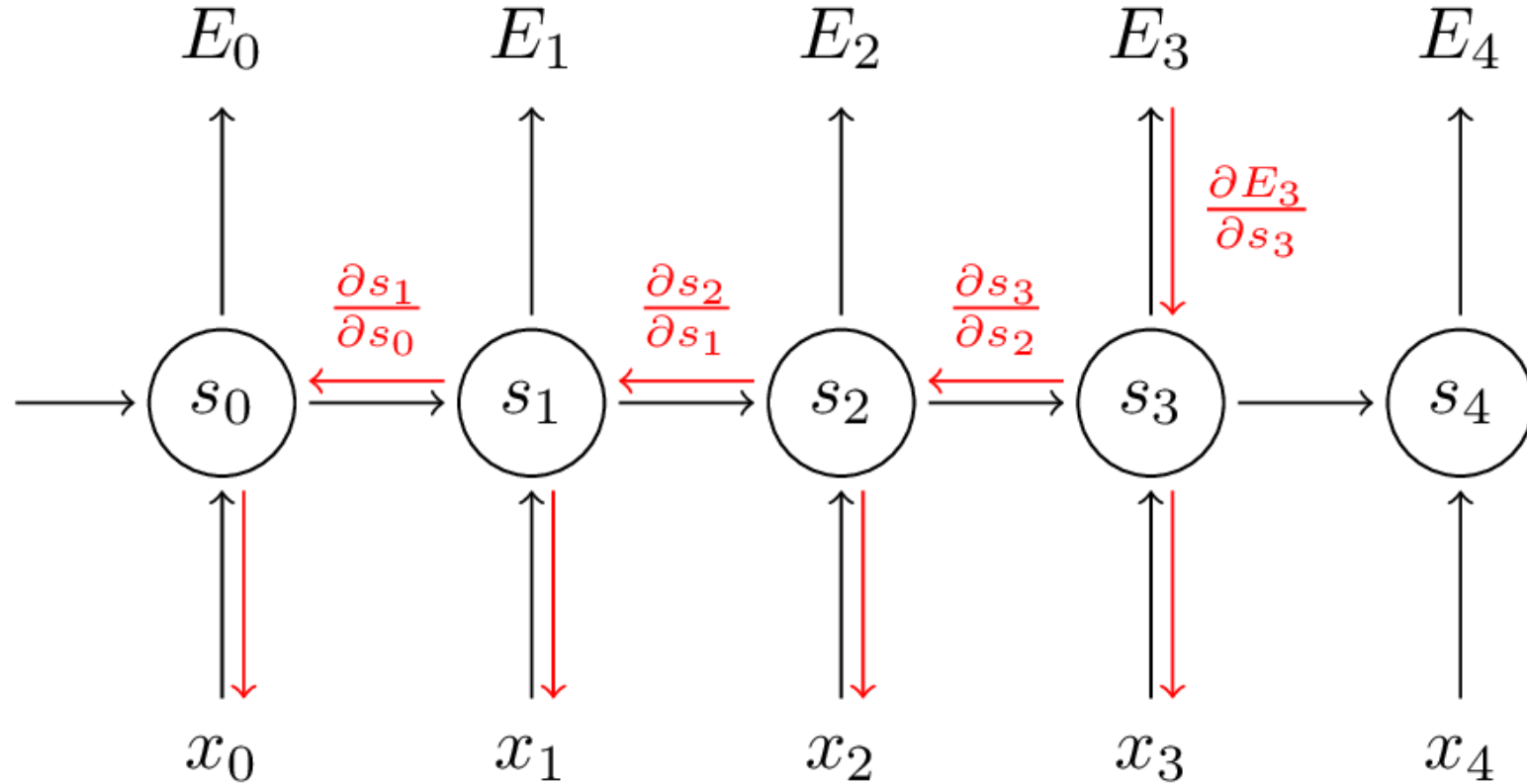


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

Update h_{t-1} to get h_t , a filtered version of the cell state C_t

Backpropagation through time (BTT)




Further possibilities

Bidirectional LSTMs
Stacked LSTMs
ConvLSTMs
Attention
Transformers

LSTMs with PyTorch


```
lstm_layer = nn.LSTM(input_dim, hidden_dim, n_layers)
```

$x_t \in \mathbb{R}^{\text{input_dim}}$
e.g. $[-3, 1, 2.5, 7, 0.17]$
input-dim = 5



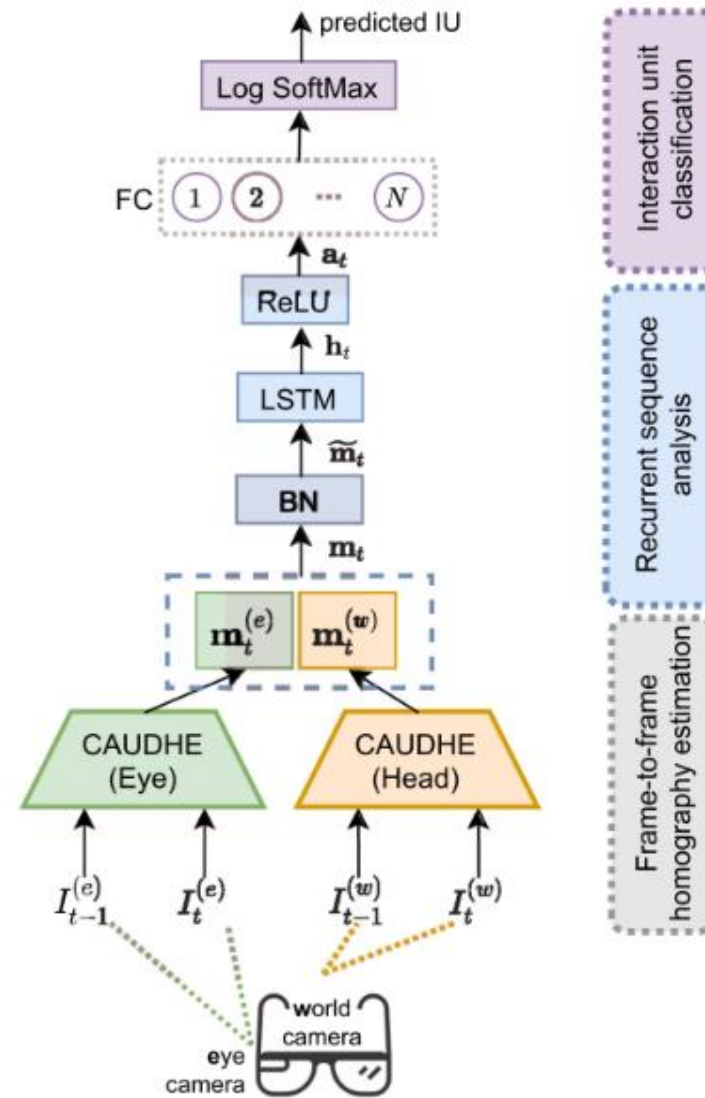
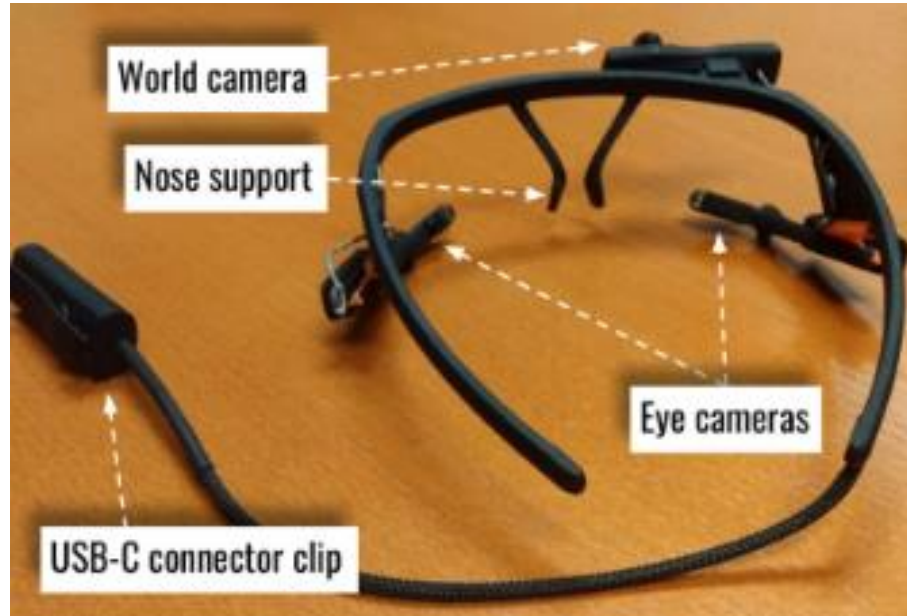
$h_t \in \mathbb{R}^{\text{hidden_dim}}$
 $c_t \in \mathbb{R}^{\text{hidden_dim}}$

$n_layers = 1$
(non-stacked LSTM)



What about the sequence length?

Example: egocentric gesture recognition



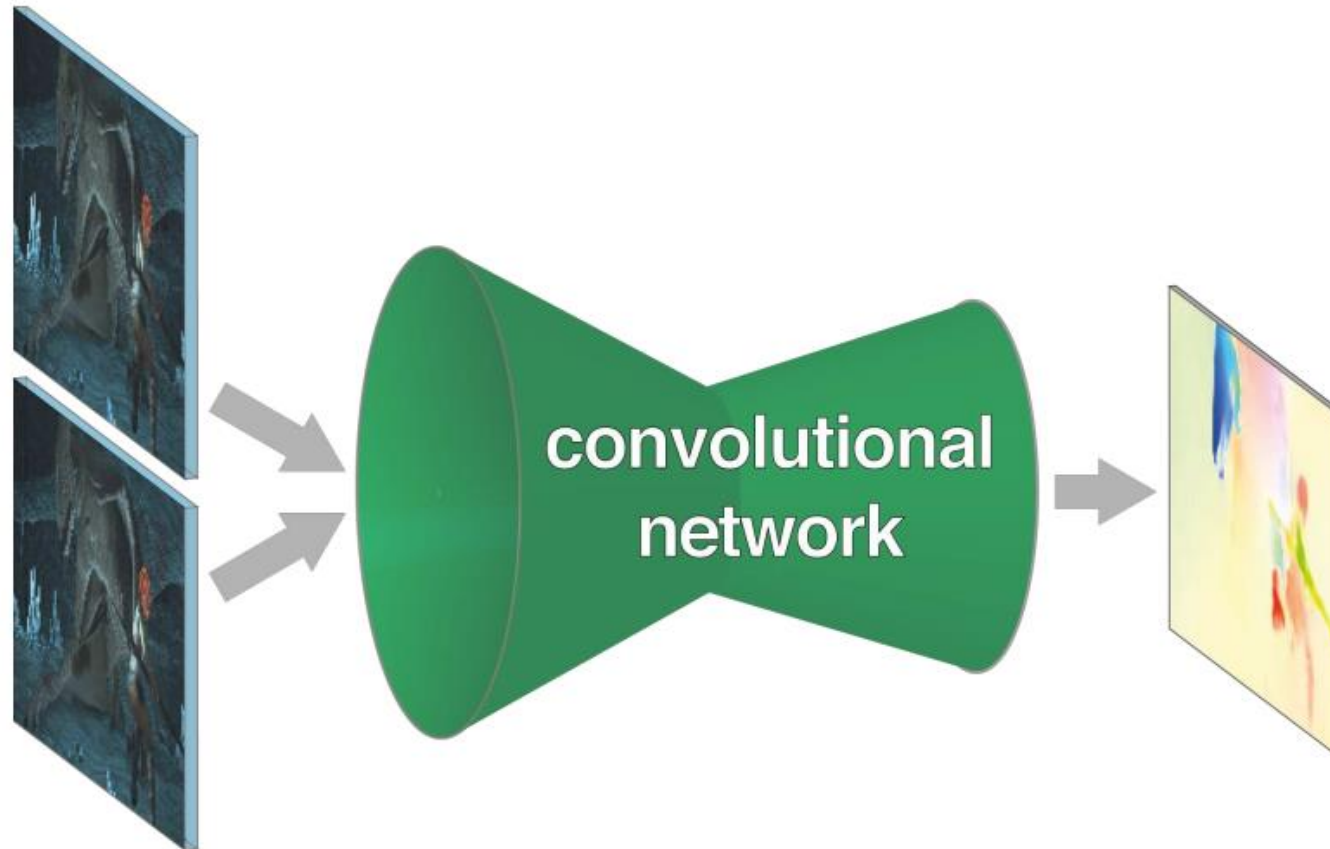
[Results \(video\)](#)

Optical flow:

FlowNet, SpyNet, LikeNet

FlowNet (ICCV2015)

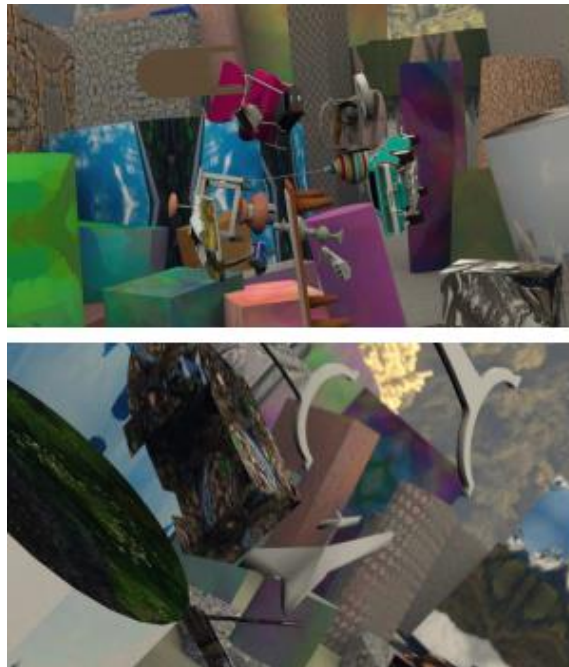
Can optic flow estimation be *learned*?



Supervised learning: what about the *labelled* data?



FlyingChairs

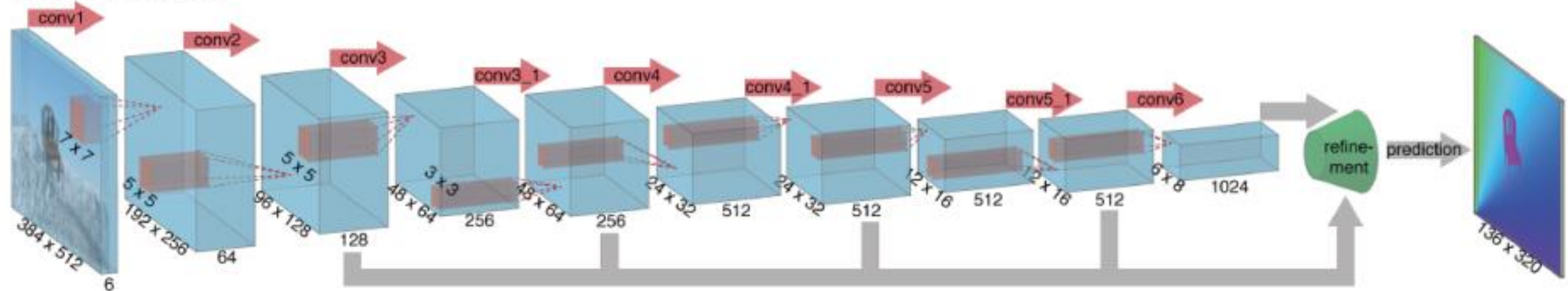


FlyingThings3D

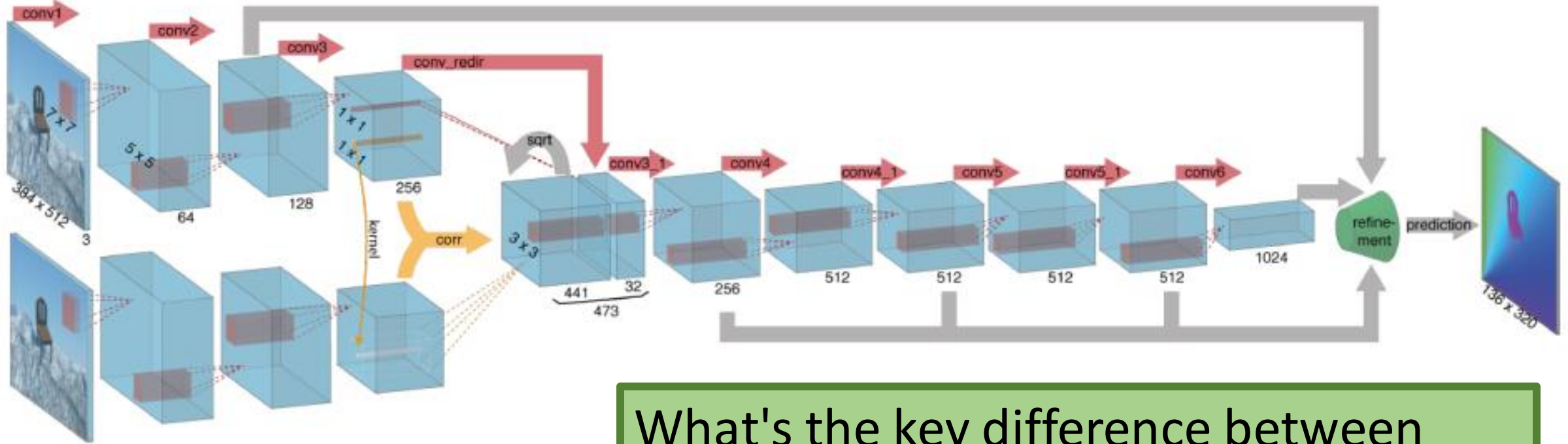


Sintel

FlowNetSimple



FlowNetCorr



What's the key difference between FlowNetSimple and FlowNetCorr?

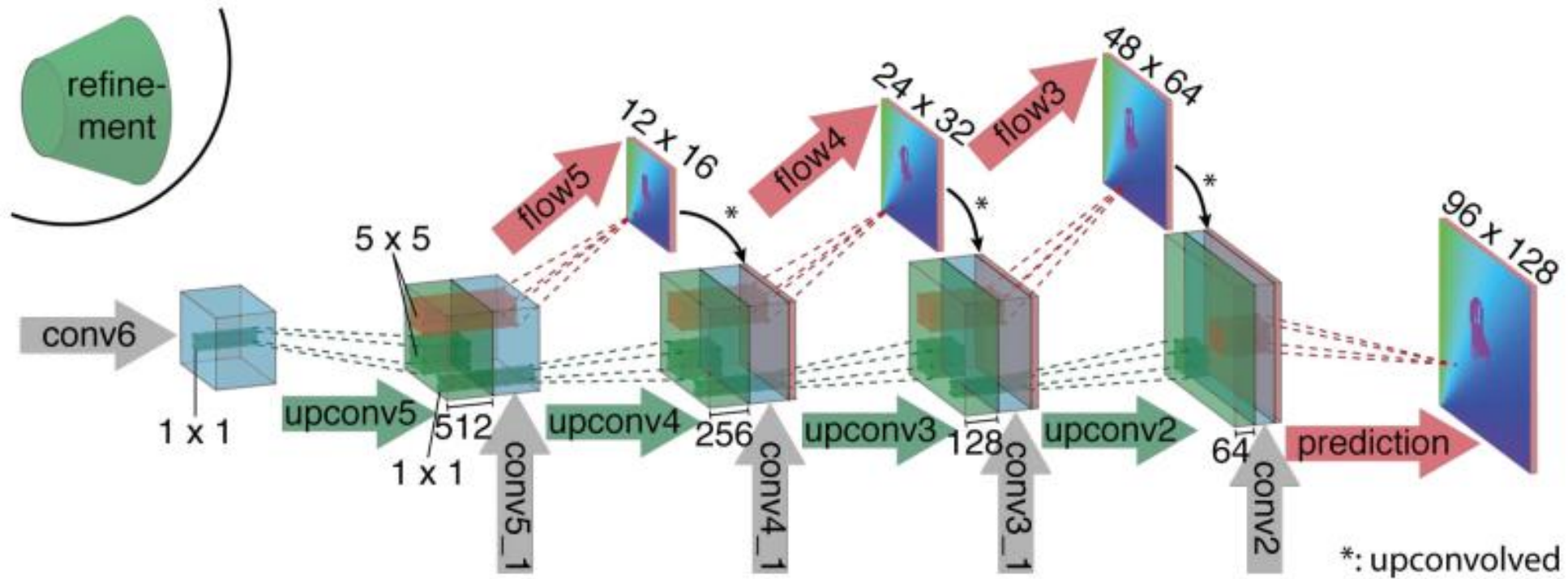
Correlation layer

$$c(\mathbf{x}_1, \mathbf{x}_2) = \sum_{\mathbf{o} \in [-k, k] \times [-k, k]} \langle \mathbf{f}_1(\mathbf{x}_1 + \mathbf{o}), \mathbf{f}_2(\mathbf{x}_2 + \mathbf{o}) \rangle$$

Like convolution between feature maps (no learnable weights)

Refinement

Upsampling



Findings

It is **possible to learn** to estimate optic flow

Training data need **not** be **realistic**

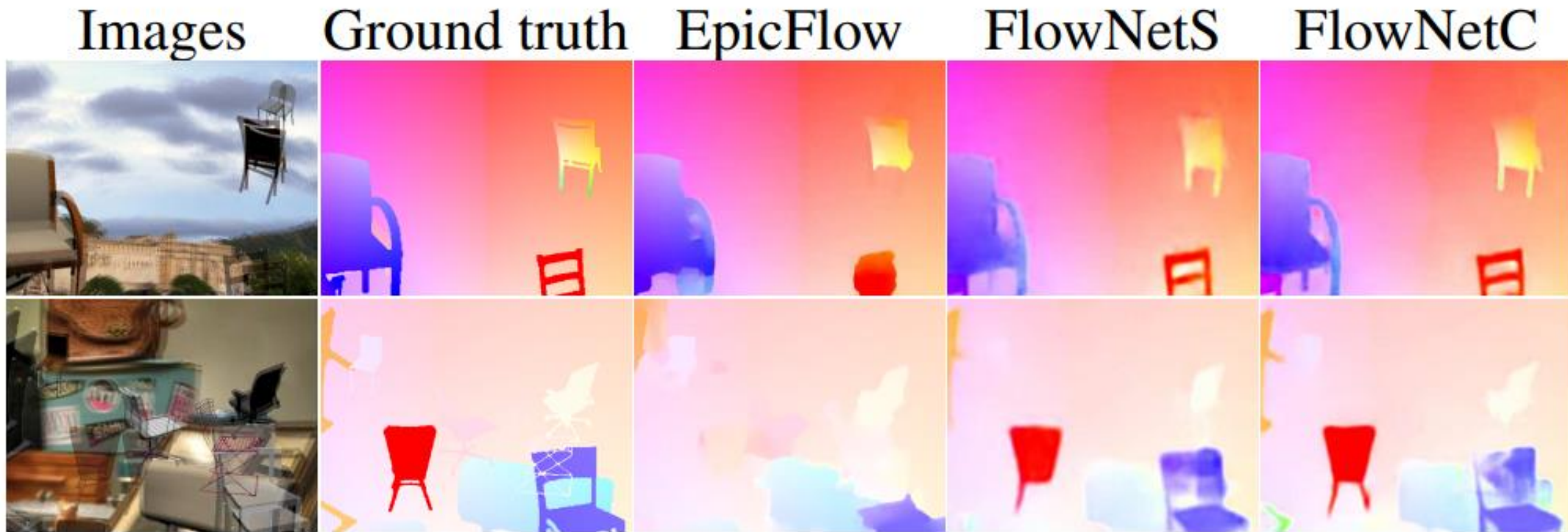
FlowNetSimple vs FlowNetCorr? --> Depends on dataset

FlowNetCorr:

- slightly overfits the training data
- has some problems with large displacements

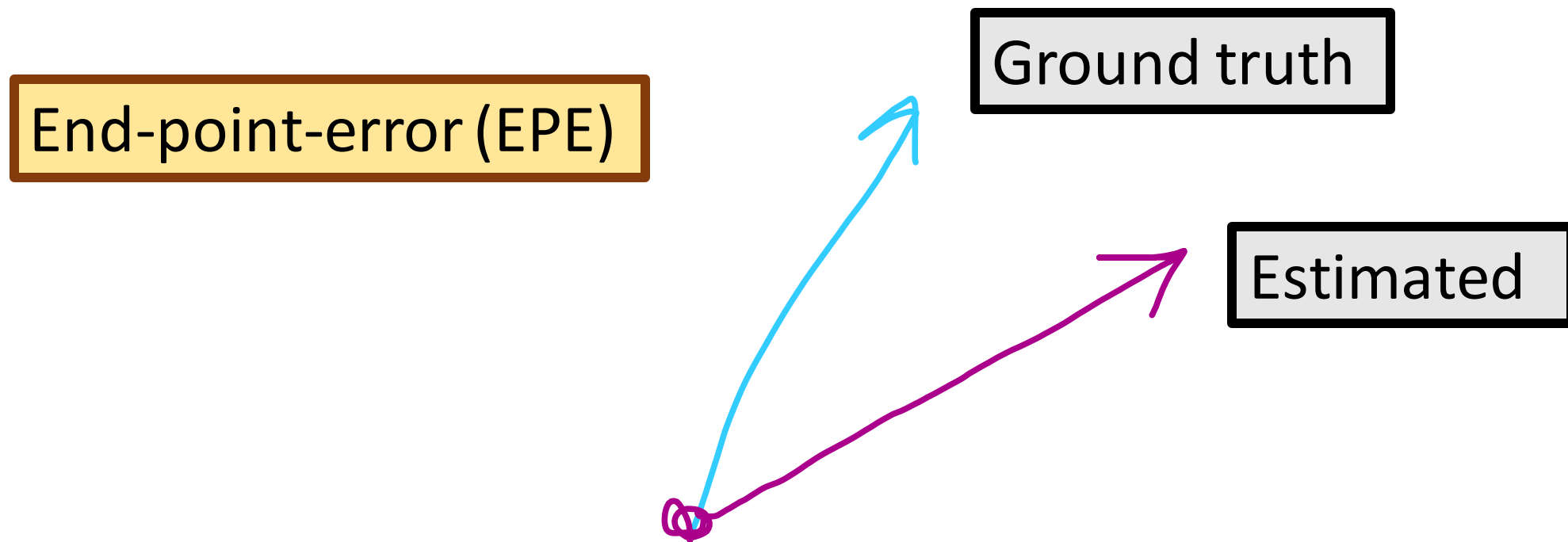
Why do you think authors relate the limitation of FlowNetCorr with large displacements with the correlation layer?

Some results



Find where FlowNetS is better than EpicFlow (CVPR 2015)
Find where FlowNetC is better than FlowNetS

Loss



FlowNet 2.0 (CVPR 2017)

FlowNet 2.0 vs FlowNet (ICCV 2015)

- schedule of presenting data
- stacked architecture with warping
- subnetwork specializing on small motions

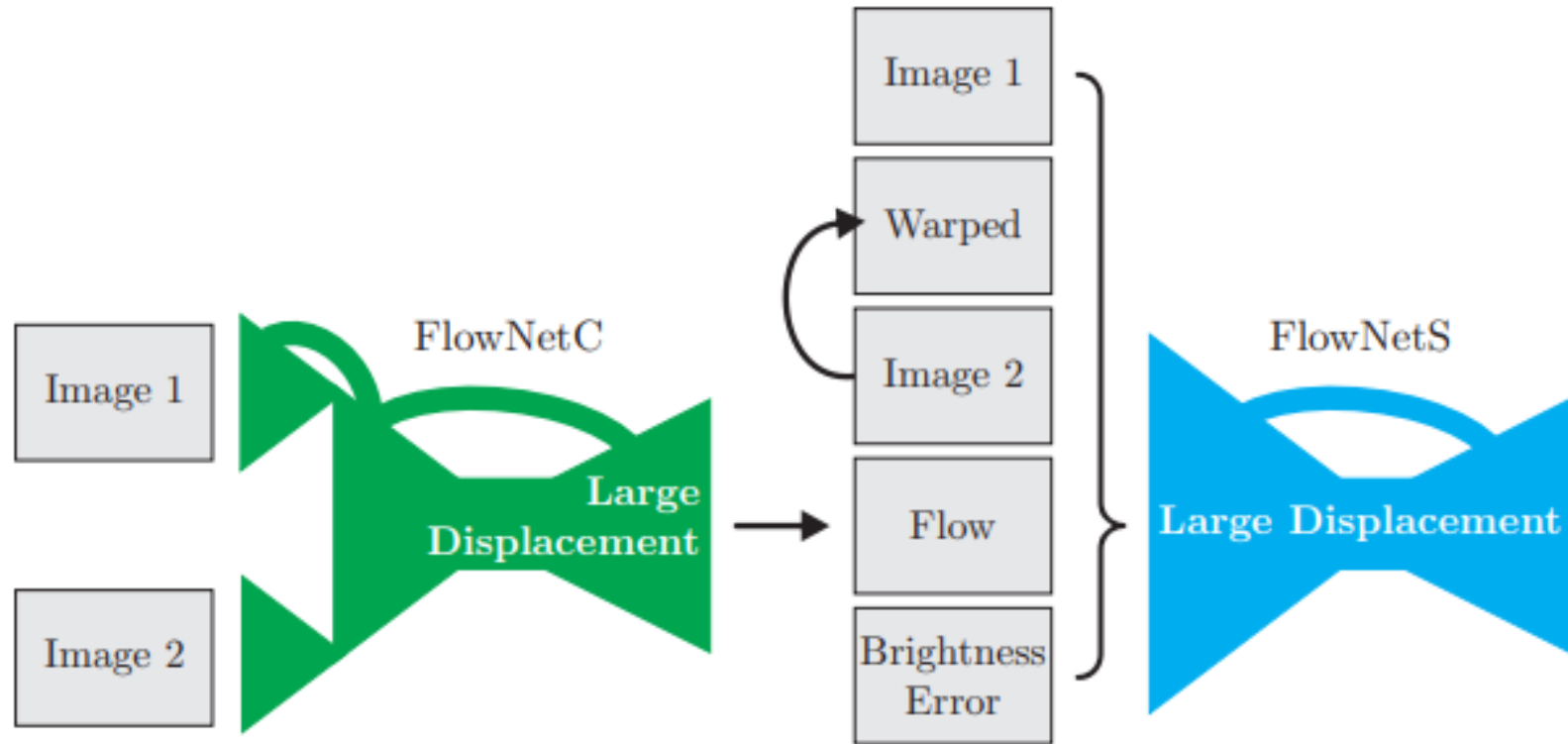
Data presentation is important

- Training only on Things3D (more realistic) --> worse results
- Best results: training on FlyingChairs first, then on Things3D

What would be a "take-home lesson" for other problems?

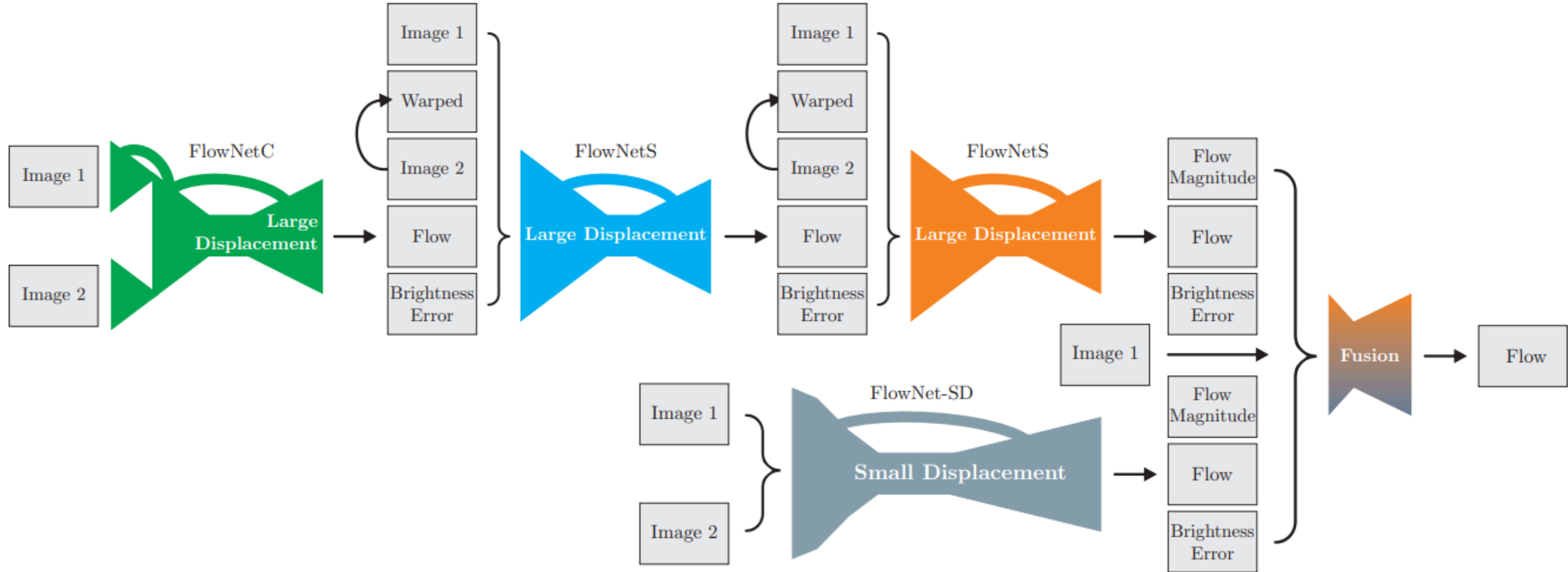
A form of curriculum learning?

Does iteration and warping help?

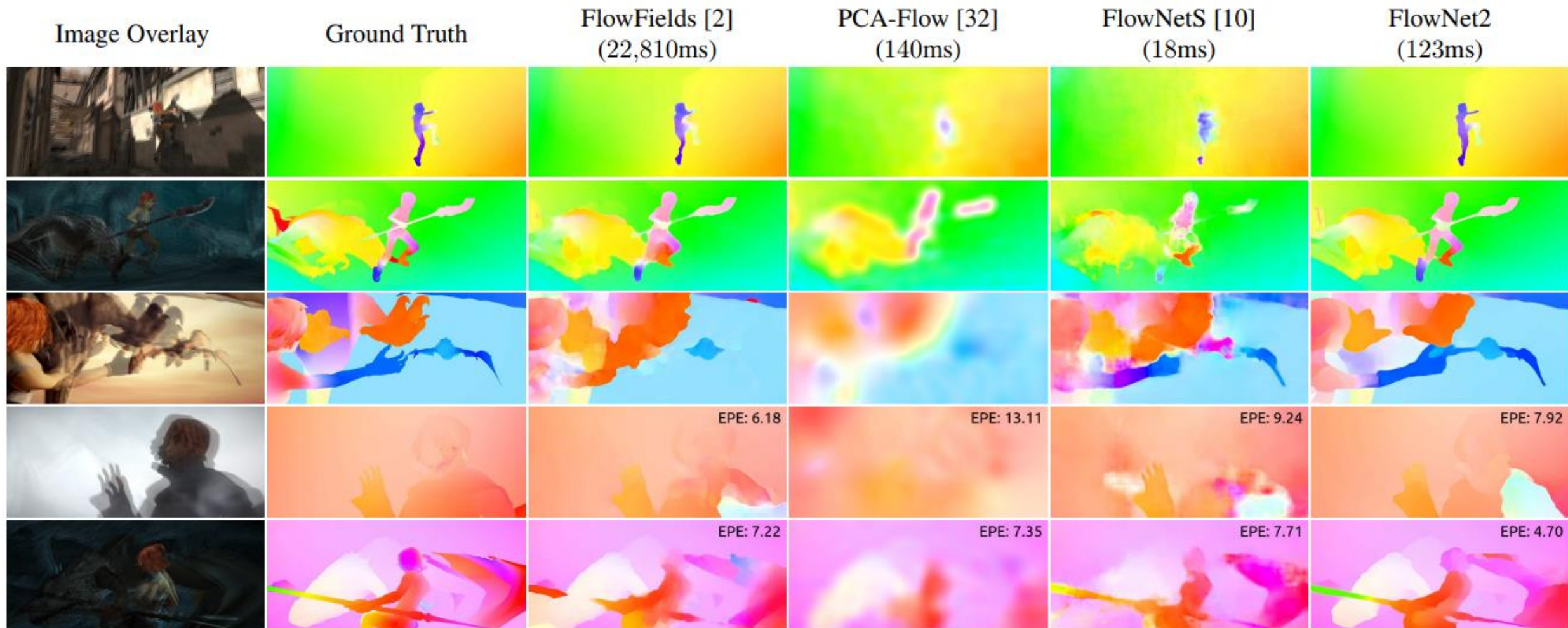


Stacking may help;
stacking+warping *always* help

Final net



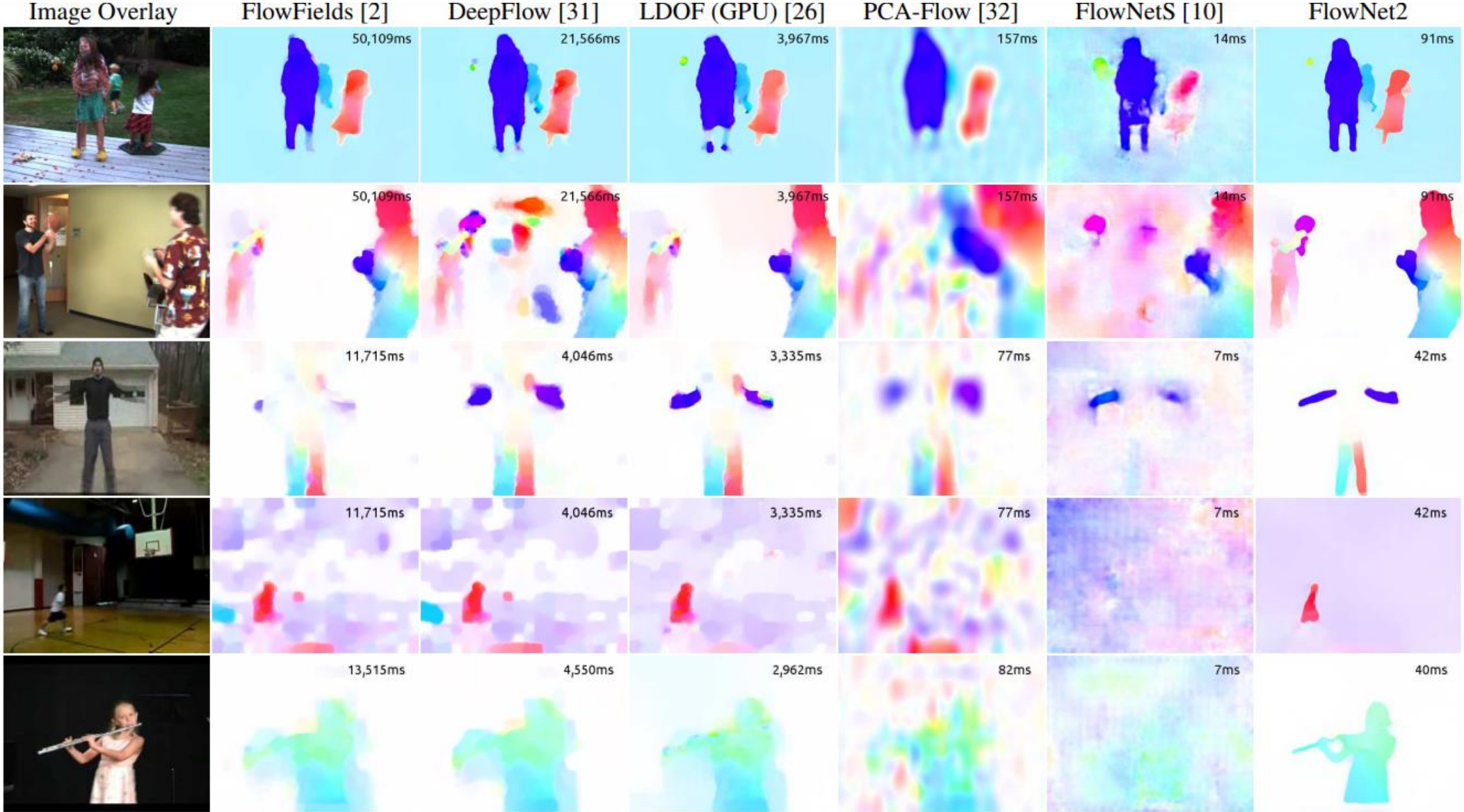
Some results on Sintel



Some results on real data

Middlebury

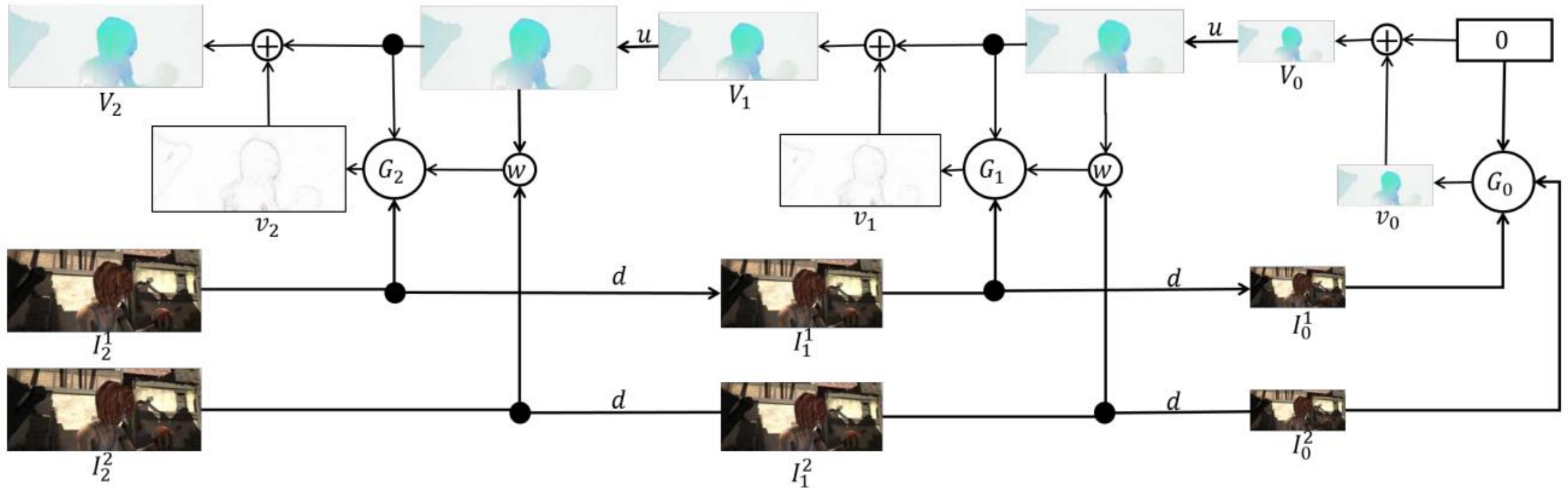
UCF101



Conclusions

- FlowNet 2.0 is marginally slower than FlowNet
- Estimation error reduced by more than 50%
- Performs on par with state-of-the-art methods
- Runs at interactive frame rates
- Faster variants run up to 140fps

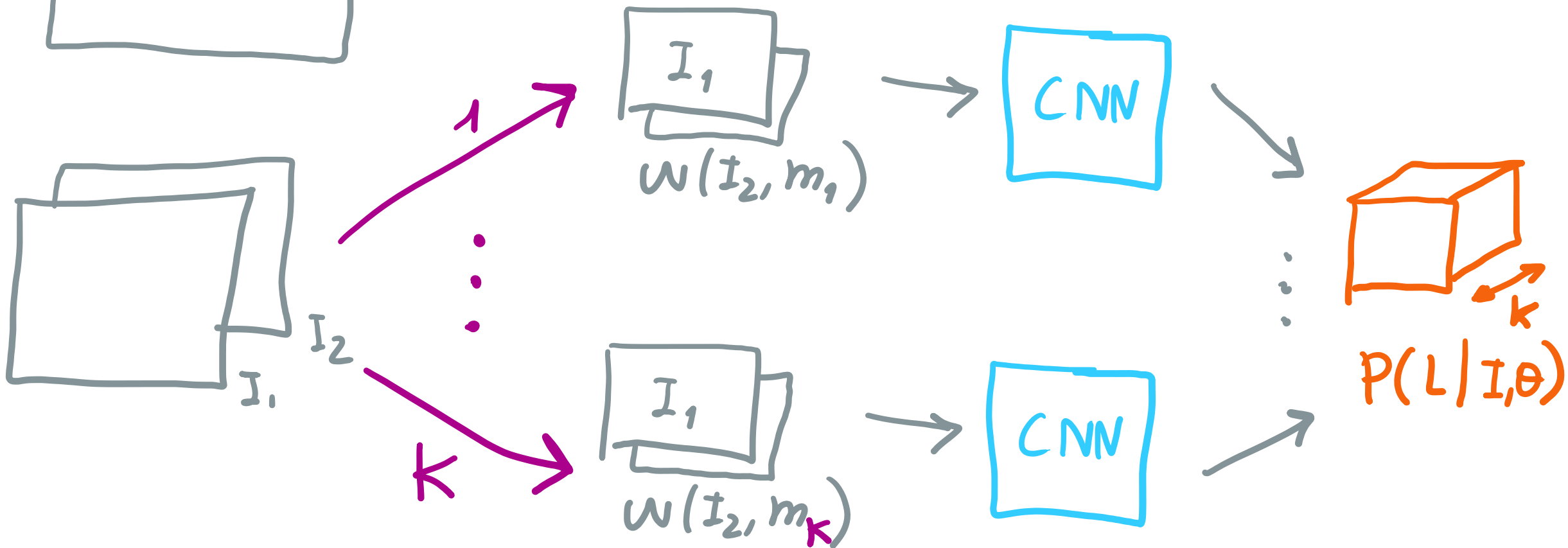
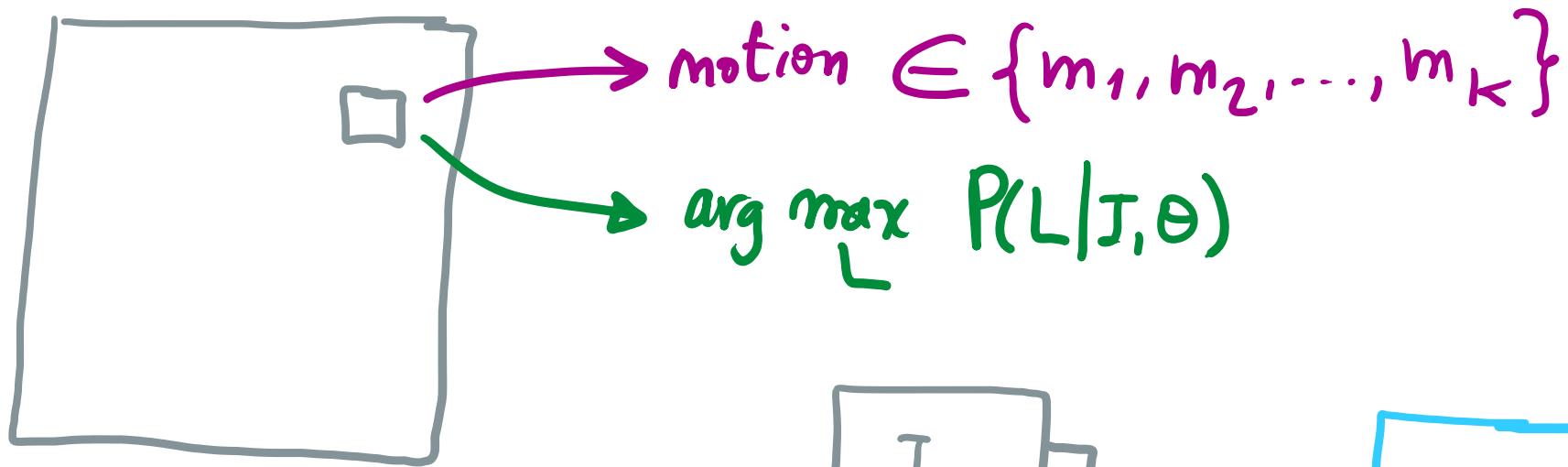
SPyNet (CVPR 2017)



- Each G_i trained independently (previous G_i already trained)
- Since each G_i assumes small motion (simpler task), less overall parameters are required

LikeNet (BMVC 2018)

- Motion as dense *classification* problem
- **Unsupervised** (no ground truth required!)
- **Siamese** architecture



Loss function?

We don't have ground truth to compare with!

$$C(I; \theta) = \sum_i \sum_{\mathbf{m}_k \in M} P(L_i = \mathbf{m}_k | I, \theta) D(i, \mathbf{m}_k)$$

$$D(i, \mathbf{m}_k) = JSD(\mathbf{F}(x_i + \mathbf{m}_k, t + dt) || \mathbf{F}(x_i, t))$$

$D(\cdot)$ could be a function of image **values**,
but with **features** works better. Any idea **why**?

Motion should be quantised!

How many branches would be required for 50 values in t_x and t_y each?

Multiscale approach: # branches: 121, 169, 49, 9, 9

More branches at lowest resolution

Comparison

Method	Number of learned parameters
FlowNetS	32,070,472
FlowNetC	32,561,032
SpyNet	1,200,250
LikeNet	697,028

- performs **better** than the *other unsupervised* methods
- **generalizes** well to unknown datasets (without finetuning)