

Videos

Saurabh Gupta

CS 543 / ECE 549 Computer Vision

Spring 2021

Outline

- Correspondence Problems
 - Optical Flow
 - Tracking
 - Mid-level Correspondence
- Recognition in Videos
- Videos as a source of supervision

Correspondence in Time

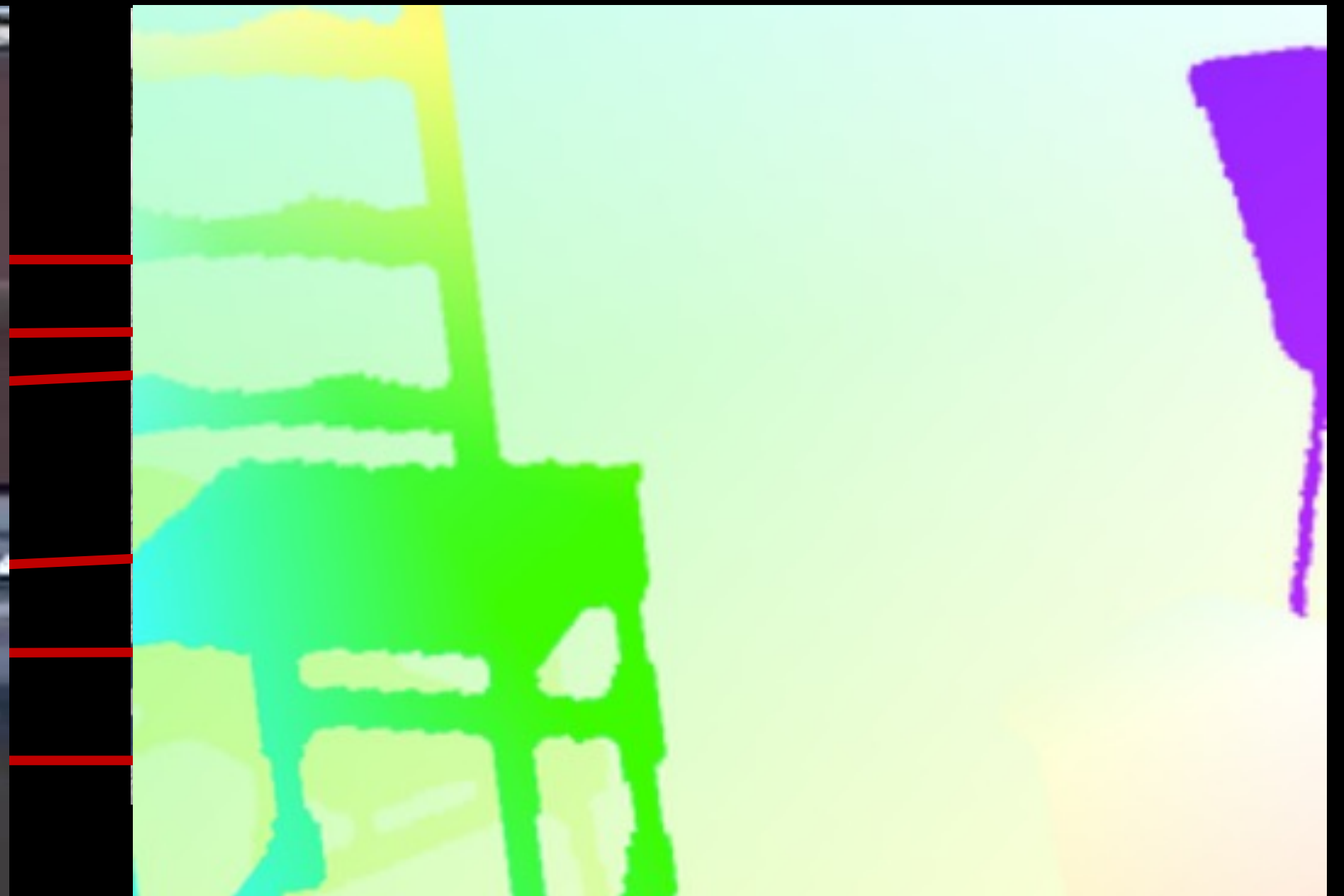
Tracking
(Box-level, long-range)

Middle Ground
(Mid-level, long-range)

Optical Flow
(Pixel-level, short-range)



Human Annotations Self-Supervised / Unsupervised Learning Synthetic Data



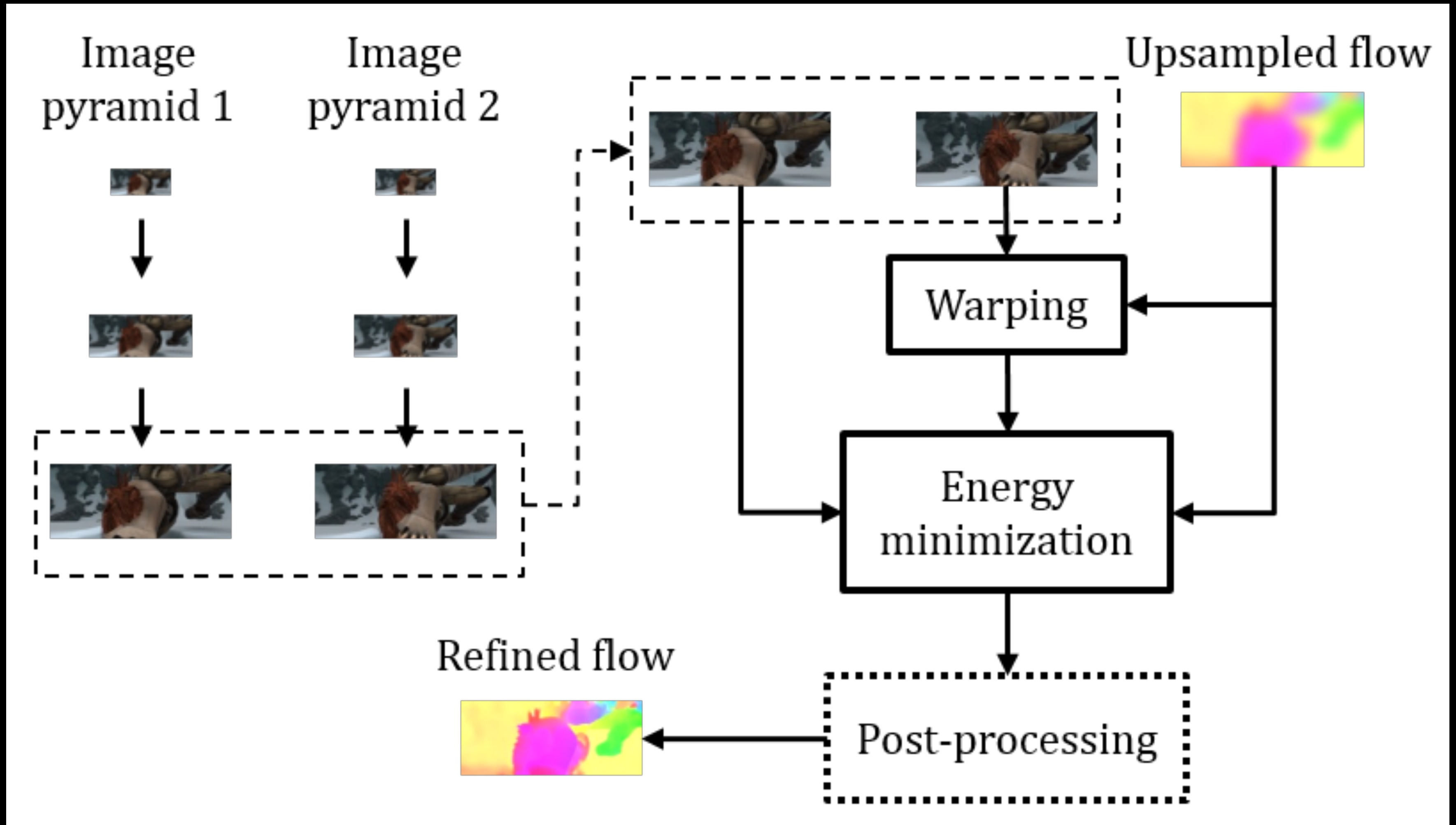
Optical Flow

- Data / Supervision
- Architecture

Datasets

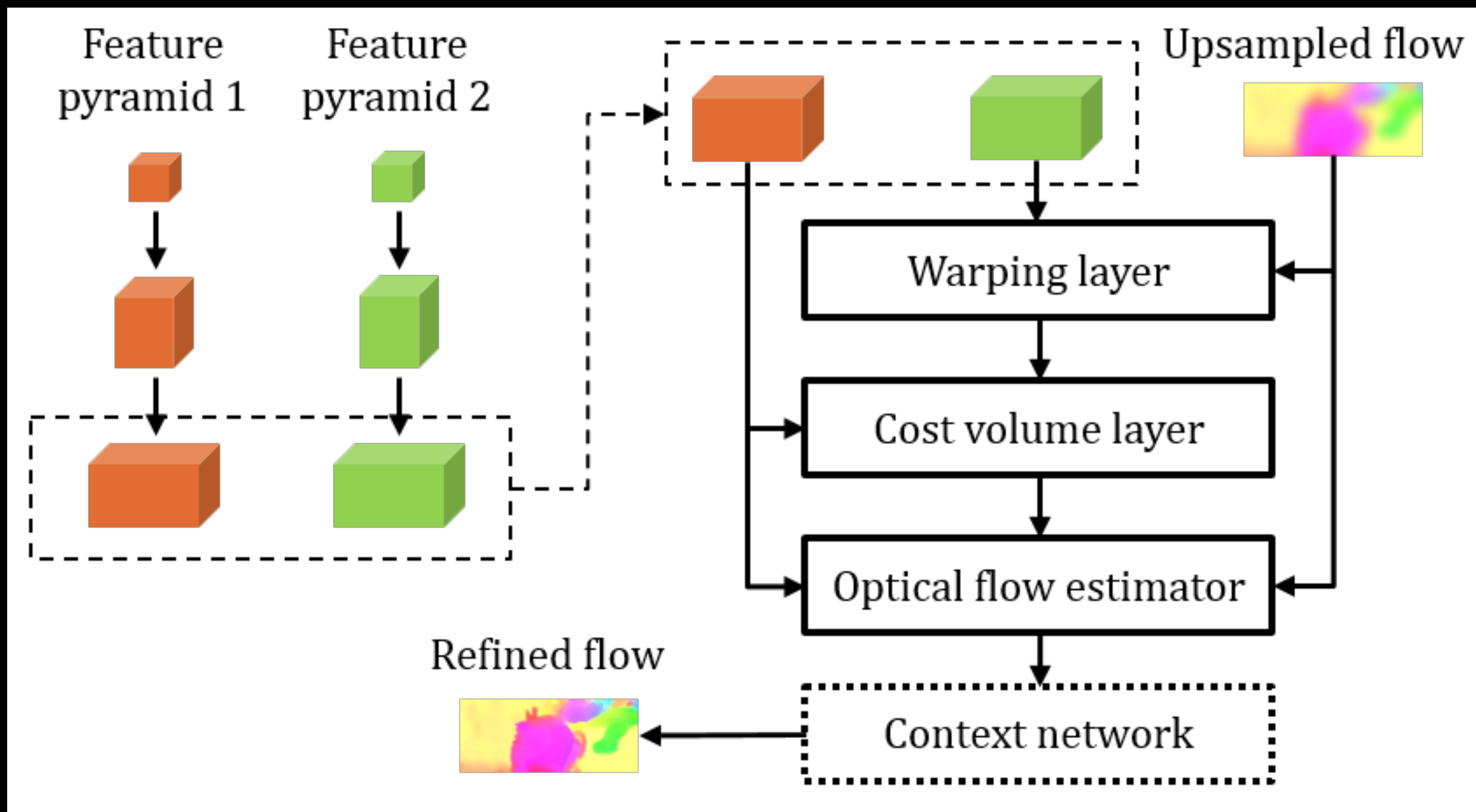
- Traditional datasets: Yosemite, Middlebury
- KITTI:
http://www.cvlibs.net/datasets/kitti/eval_scene_flow.php?benchmark=flow
- Sintel: <http://sintel.is.tue.mpg.de/>
- Synthetic Datasets
 - Flying Chairs et al: <https://lmb.informatik.uni-freiburg.de/resources/datasets/FlyingChairs.en.html>
- Supervision: from Simulation
- Metrics: End-point Error

“Classical Optical Flow Pipeline”



PWC Net

$$cv^l(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{N} \left(\mathbf{c}_1^l(\mathbf{x}_1) \right)^\top \mathbf{c}_w^l(\mathbf{x}_2),$$



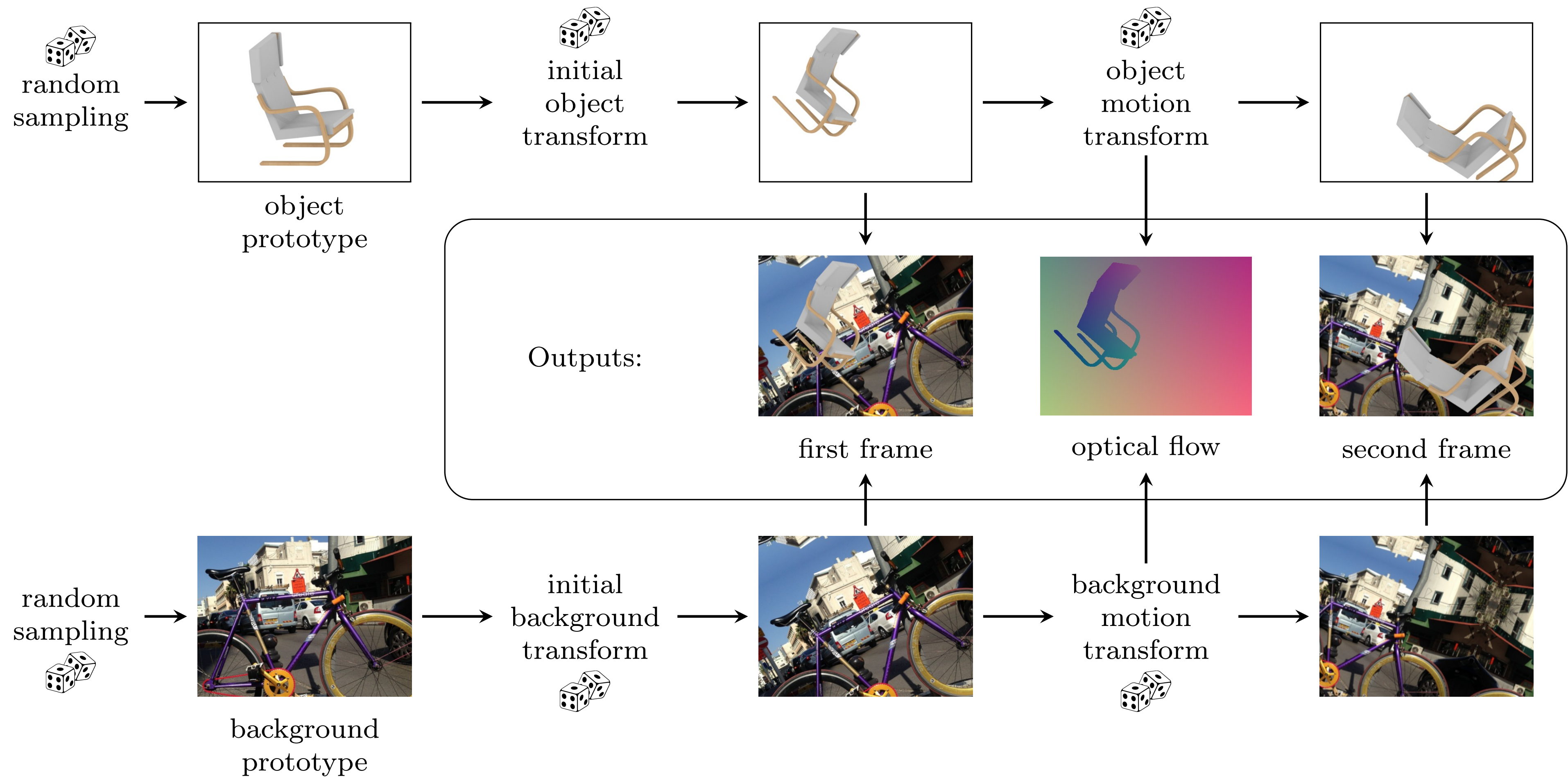
PWC Net



Max. Disp.	Chairs	Sintel Clean	Sintel Final	KITTI 2012 AEPE	KITTI 2012 Fl-all	KITTI 2015 AEPE	KITTI 2015 Fl-all
0	2.13	3.66	5.09	5.25	29.82%	13.85	43.52%
2	2.09	3.30	4.50	5.26	25.99%	13.67	38.99%
Full model (4)	2.00	3.33	4.59	5.14	28.67%	13.20	41.79%
6	1.97	3.31	4.60	4.96	27.05%	12.97	40.94%

(b) **Cost volume.** Removing the cost volume (0) results in moderate performance loss. PWC-Net can handle large motion using a small search range to compute the cost volume.

Flying Chairs Dataset



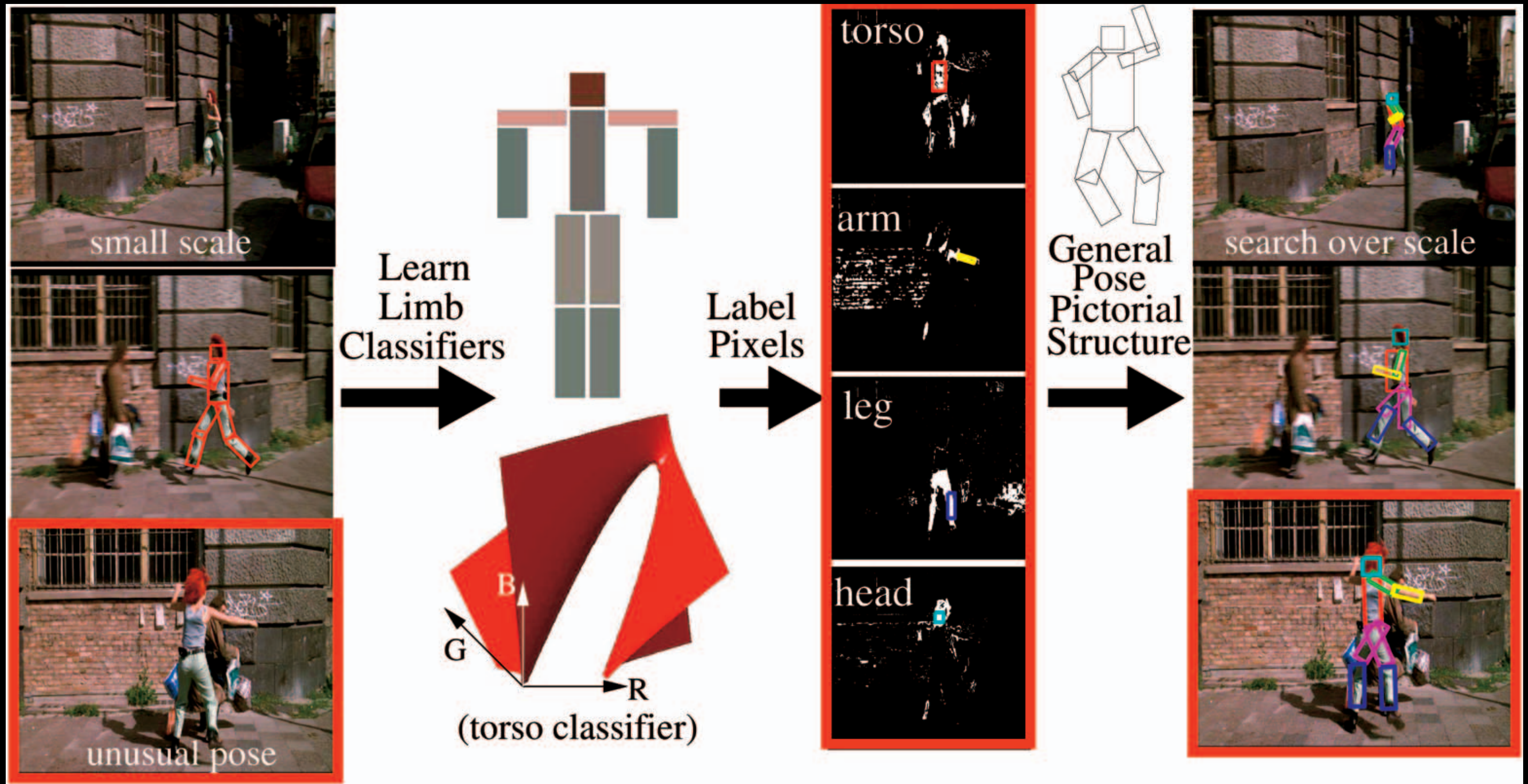
Tracking

- Problem Statements
- Tracking by Detection
- Tracking by Matching

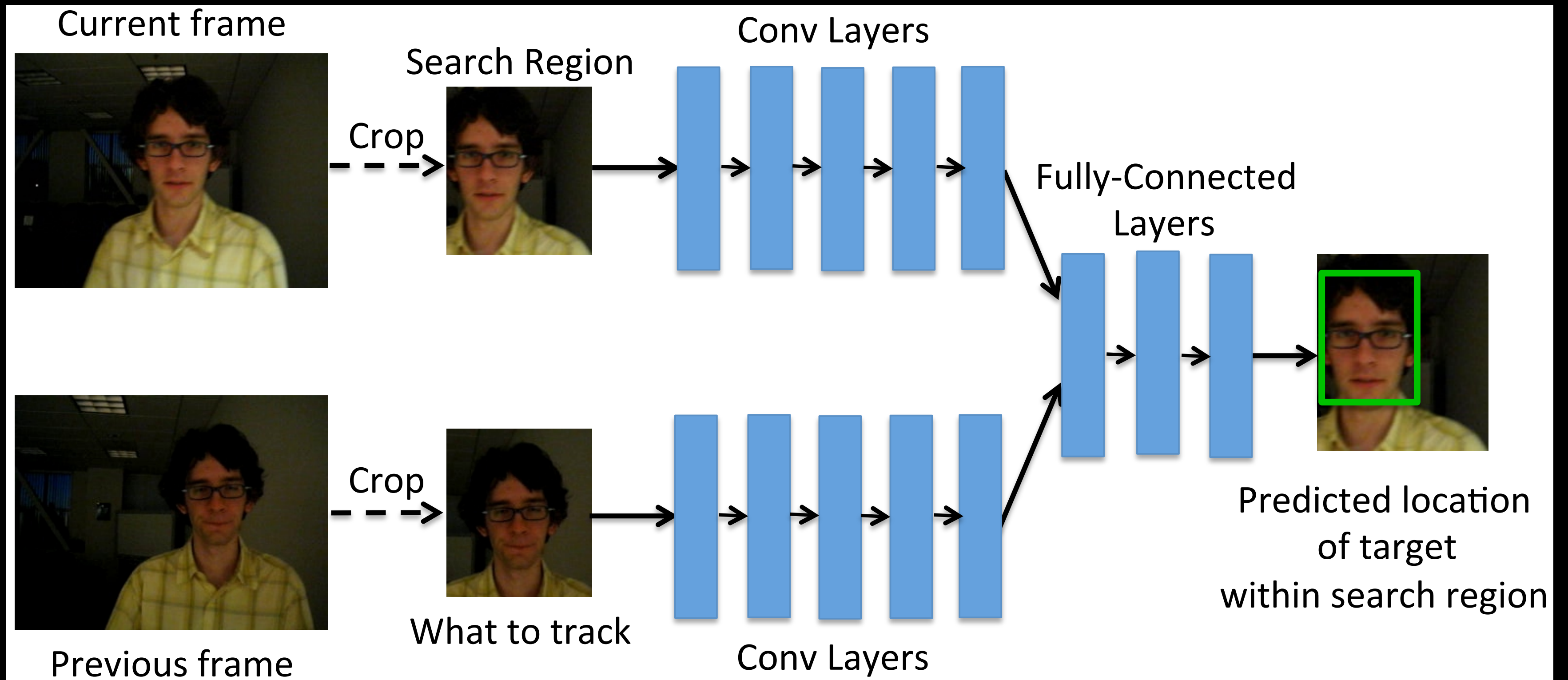
Problem Statements

- Single Object Tracking (eg: https://nanonets.com/blog/content/images/2019/07/messi_football_track.gif)
- Multi-object Tracking (eg: <https://motchallenge.net/vis/MOT20-02/gt/>)
- Multi-object Tracking and Segmentation (eg: https://www.youtube.com/watch?v=K38_pZw_P9s)

Tracking by Detection

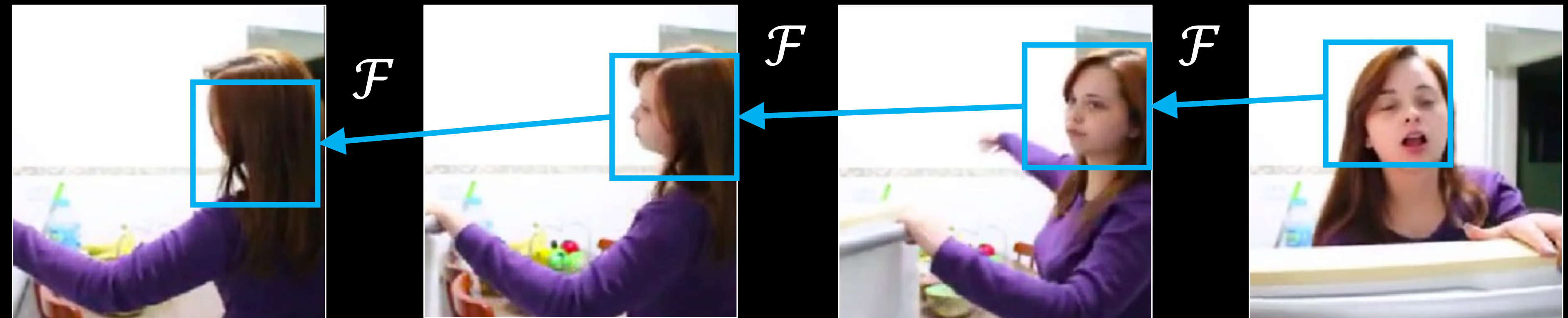


General Object Tracking



Tracking by Learning to Match

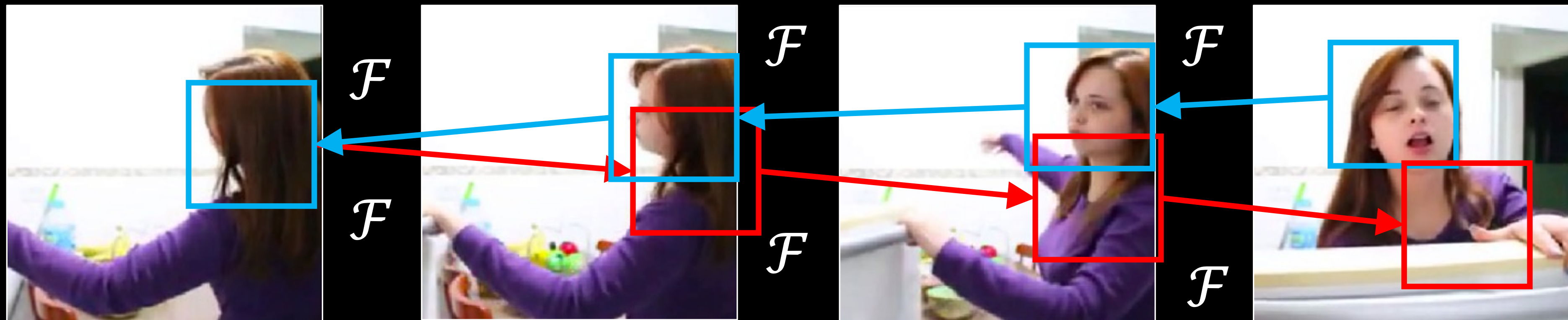
\mathcal{F} : a deep tracker



How to obtain supervision?

Supervision: Cycle-Consistency in Time

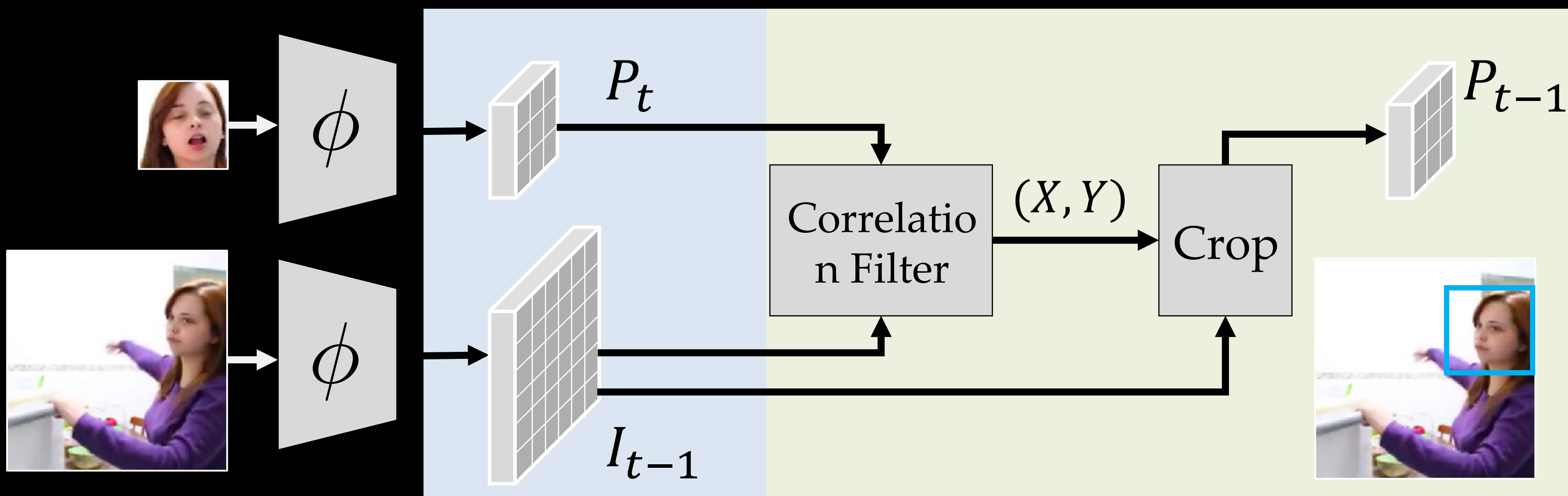
Track backwards



Track forwards, back to the future

Tracker \mathcal{F}

Densely match features in learned feature space

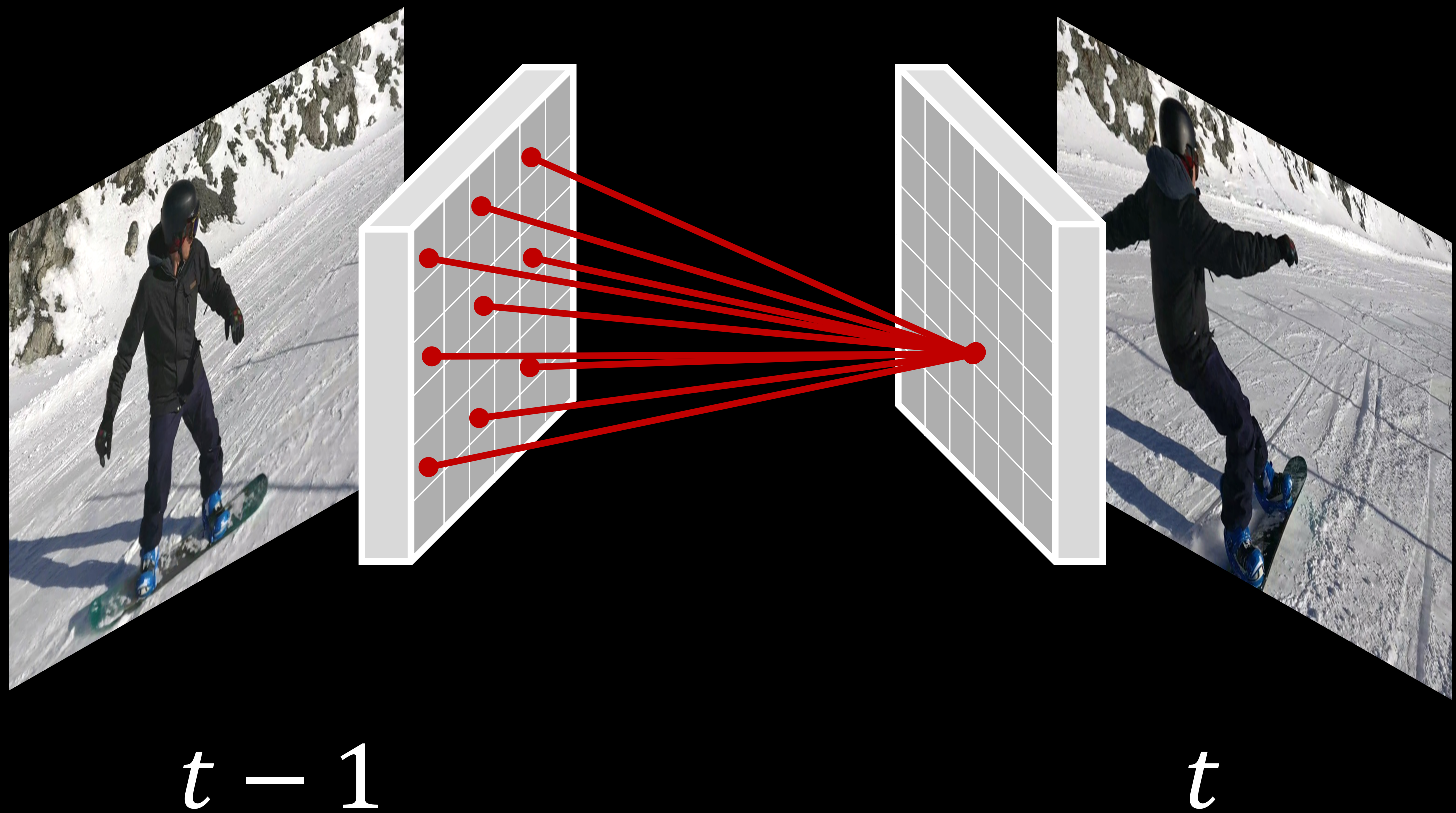


Visualization of Training

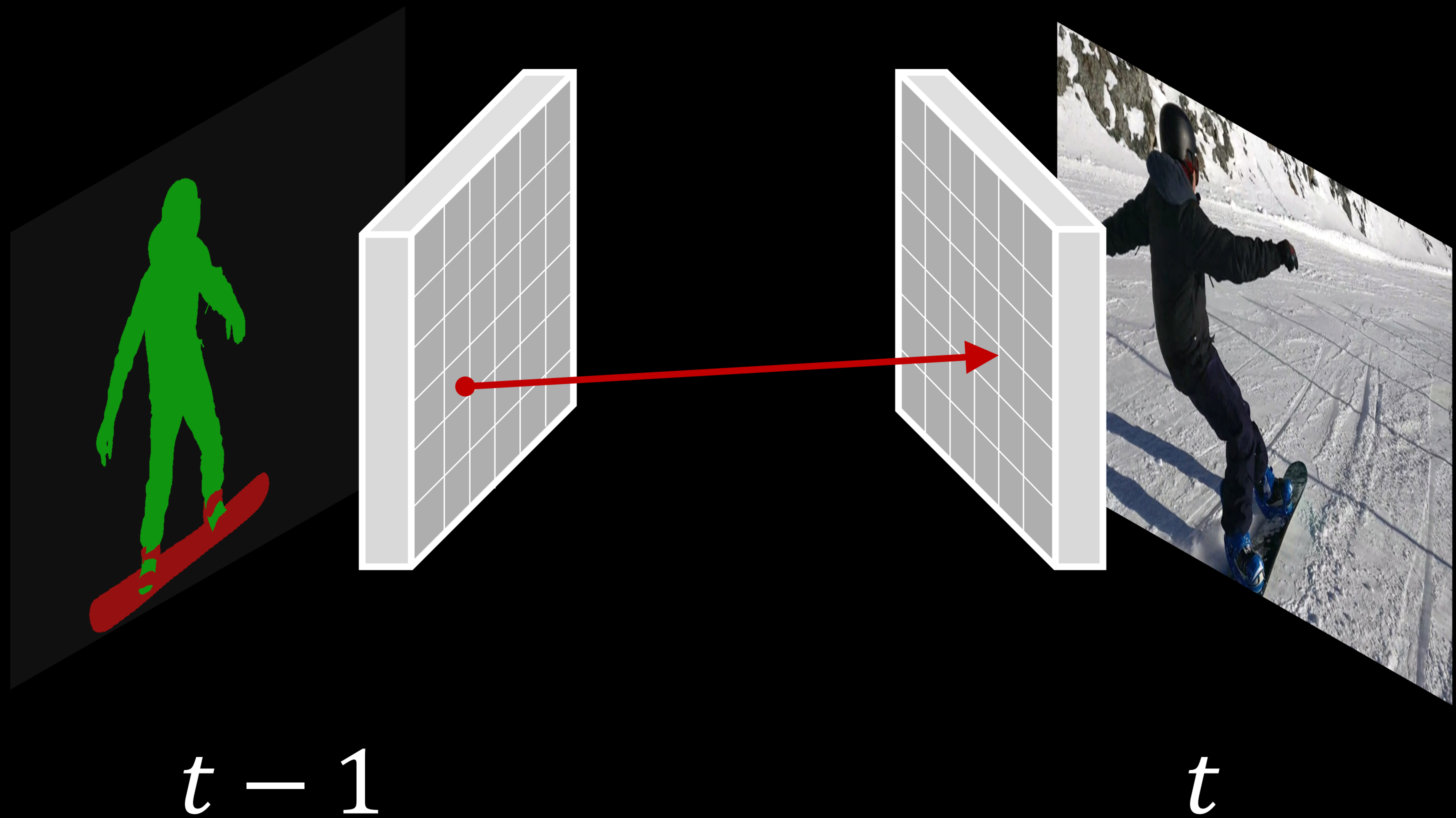


Iteration: 1200

Test Time: Nearest Neighbors in Feature Space ϕ



Test Time: Nearest Neighbors in Feature Space ϕ



Texture Tracking

DAVIS Dataset



Source: Xiaolong Wang

DAVIS Dataset: Pont-Tuset et al. *The 2017 DAVIS Challenge on Video Object Segmentation*. 2017.

Outline

- Correspondence Problems
 - Optical Flow
 - Tracking
 - Mid-level Correspondence
- Recognition in Videos
- Videos as a source of supervision

Recognition in Videos

- Tasks / Datasets
- Models

Tasks and Datasets

- Action Classification
 - Kinetics Dataset: <https://arxiv.org/pdf/1705.06950.pdf>
 - ActivityNet, Sports-8M, ...
- Action “Detection”
 - In space, in time. Eg: JHMDB, AV

Tasks and Datasets

- Time scale

- Atomic Visual Actions (AVA)

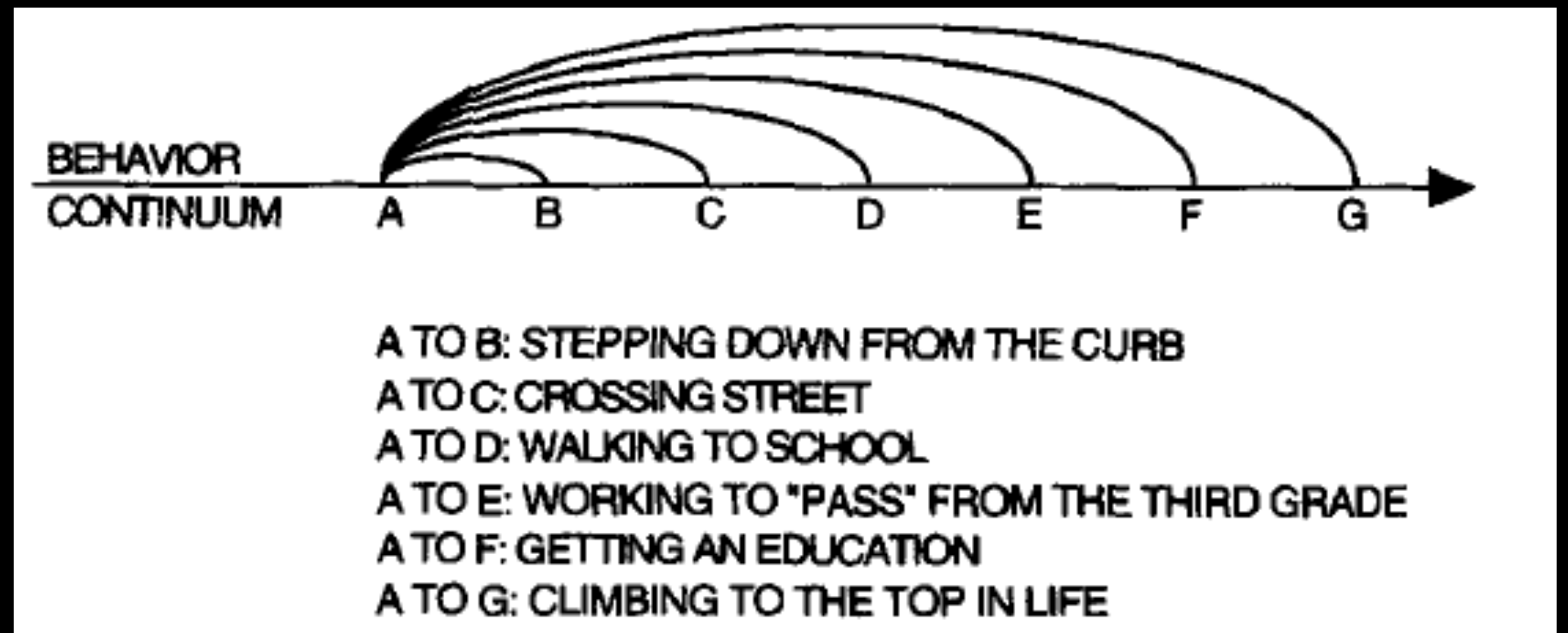
Dataset:

<https://research.google.com/ava/explore.html>

- Bias


- Something Something Dataset:

<https://20bn.com/datasets/something-something>

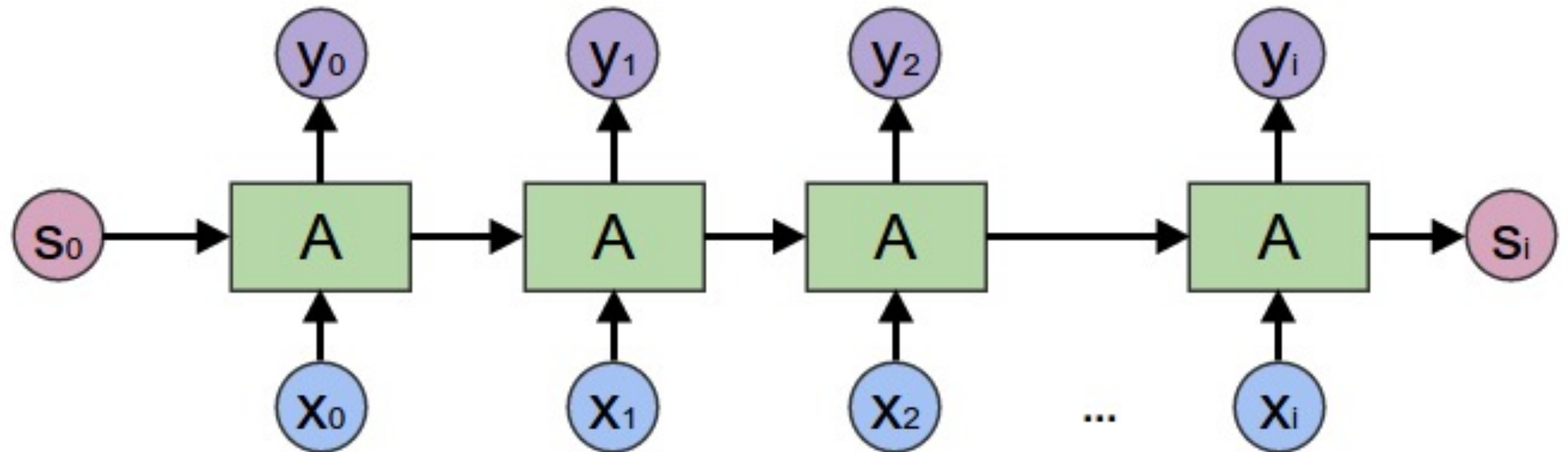


We don't quite know how to define good meaningful tasks for videos. More on this later.

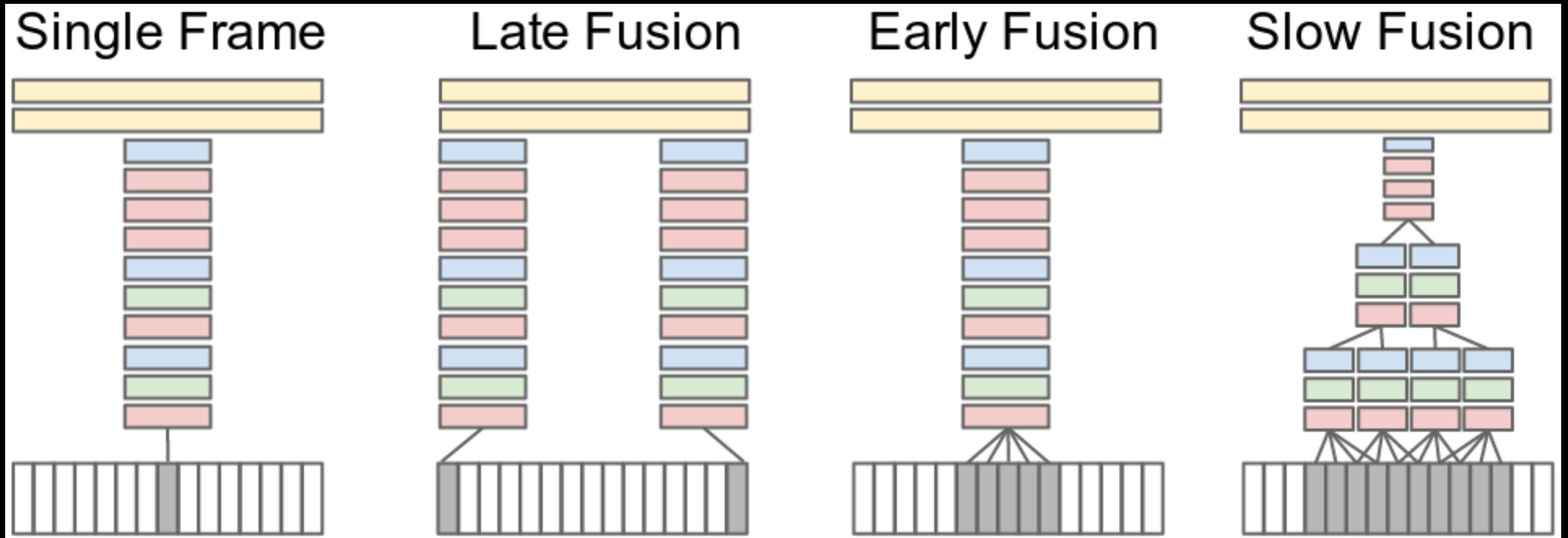
Models

- Recurrent Neural Nets (See: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>)
- Simple Extensions of 2D CNNs The diagram consists of three parts: (a) 2D convolution showing a 2D input grid, a 2D kernel, and a 2D output grid; (b) 3D convolution on multiple frames showing a 3D input volume, a 3D kernel, and a 3D output volume; (c) 3D convolution showing a 3D input volume, a 3D kernel, and a 3D output volume.
- 3D Convolution Networks
- Two-Stream Networks
- Inflated 3D Conv Nets
- Slow Fast Networks
- Non-local Networks

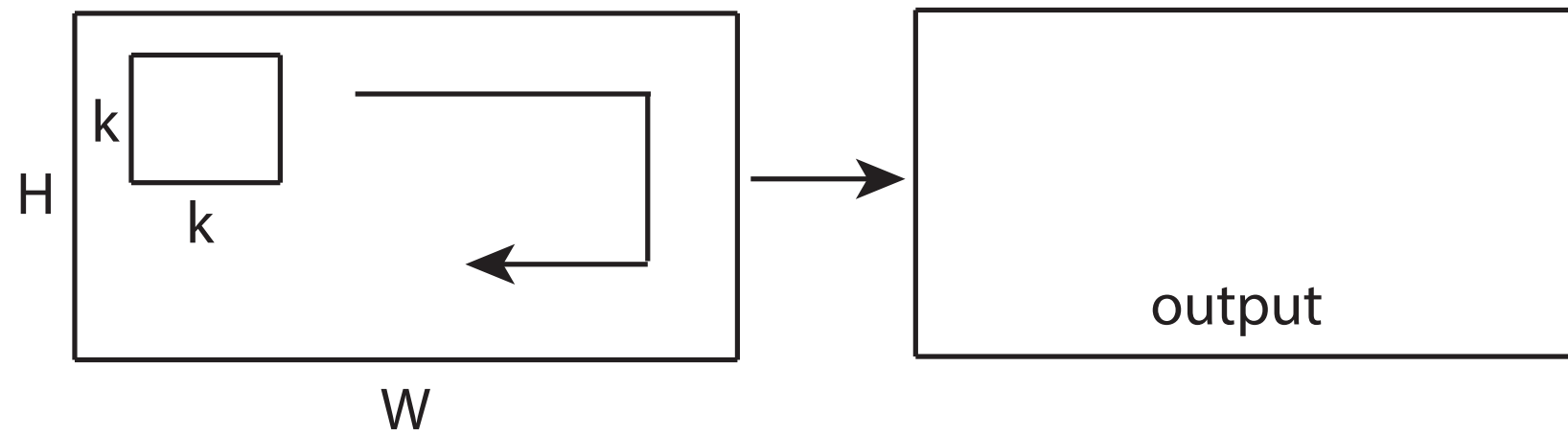
Recurrent Neural Networks



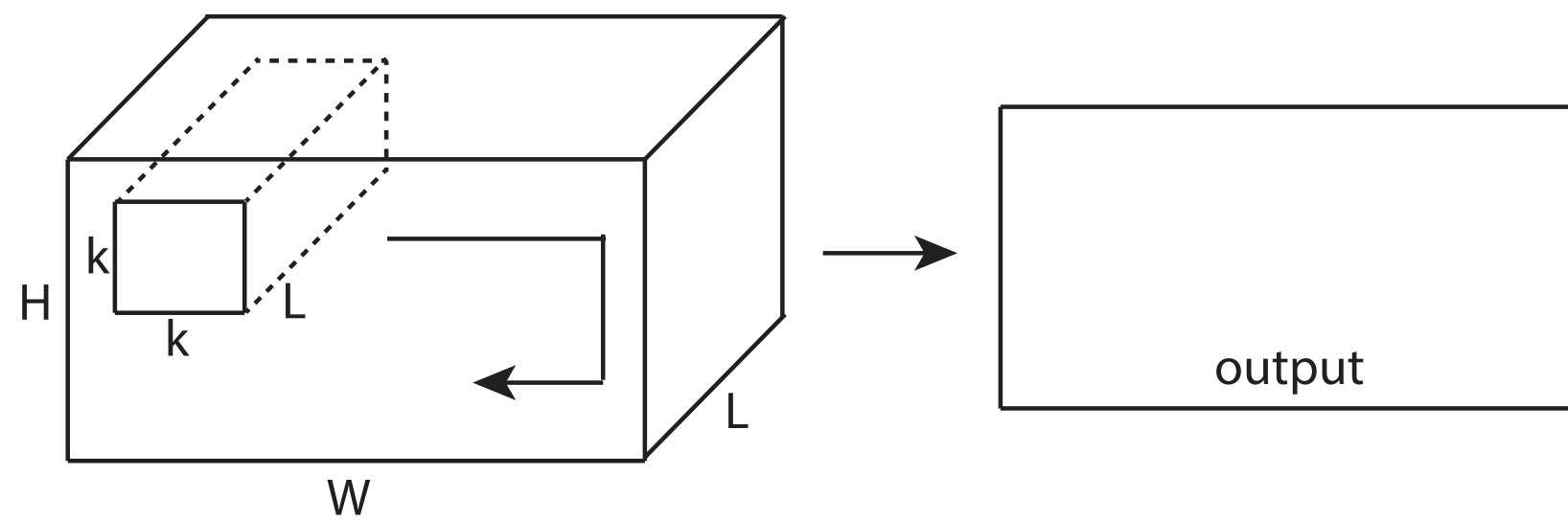
3D Convolutions



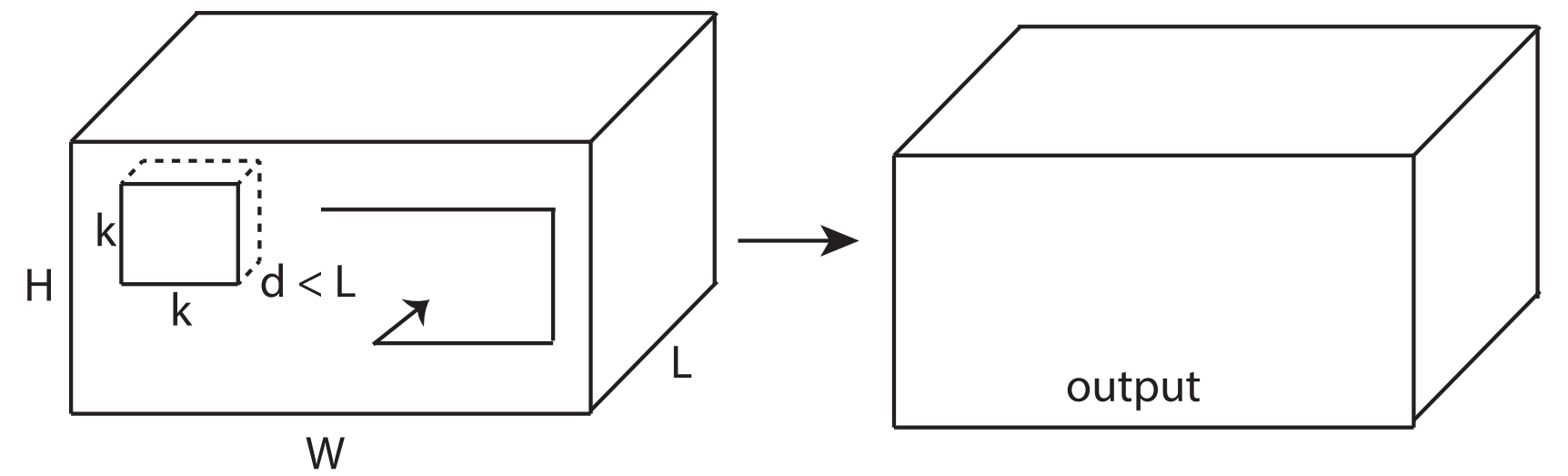
3D Convolutions



(a) 2D convolution

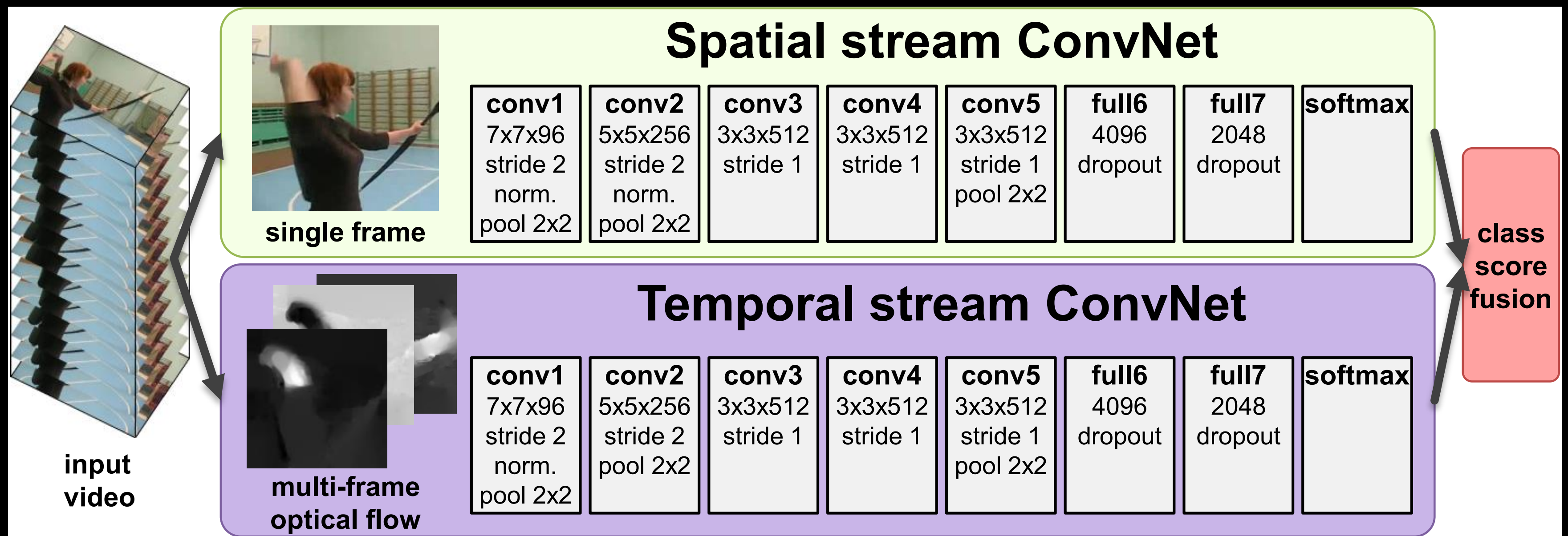


(b) 2D convolution on multiple frames



(c) 3D convolution

Two Stream Networks



Two Stream Networks

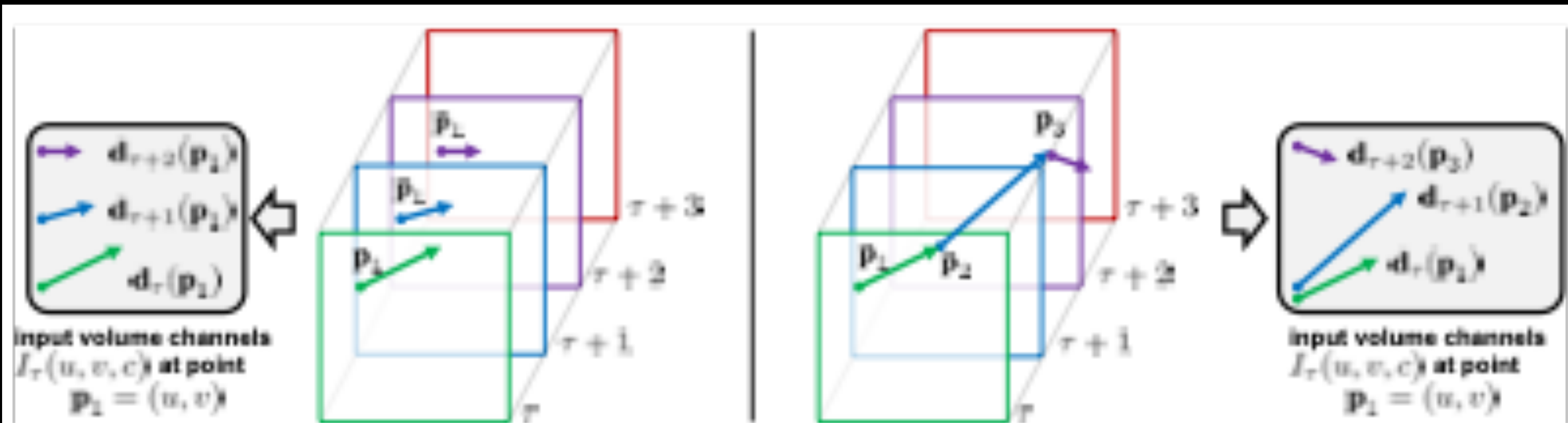


Figure 3: **ConvNet input derivation from the multi-frame optical flow.** *Left:* optical flow stacking (1) samples the displacement vectors \mathbf{d} at the same location in multiple frames. *Right:* trajectory stacking (2) samples the vectors along the trajectory. The frames and the corresponding displacement vectors are shown with the same colour.

Two Stream Networks

Table 1: Individual ConvNets accuracy on UCF-101 (split 1).

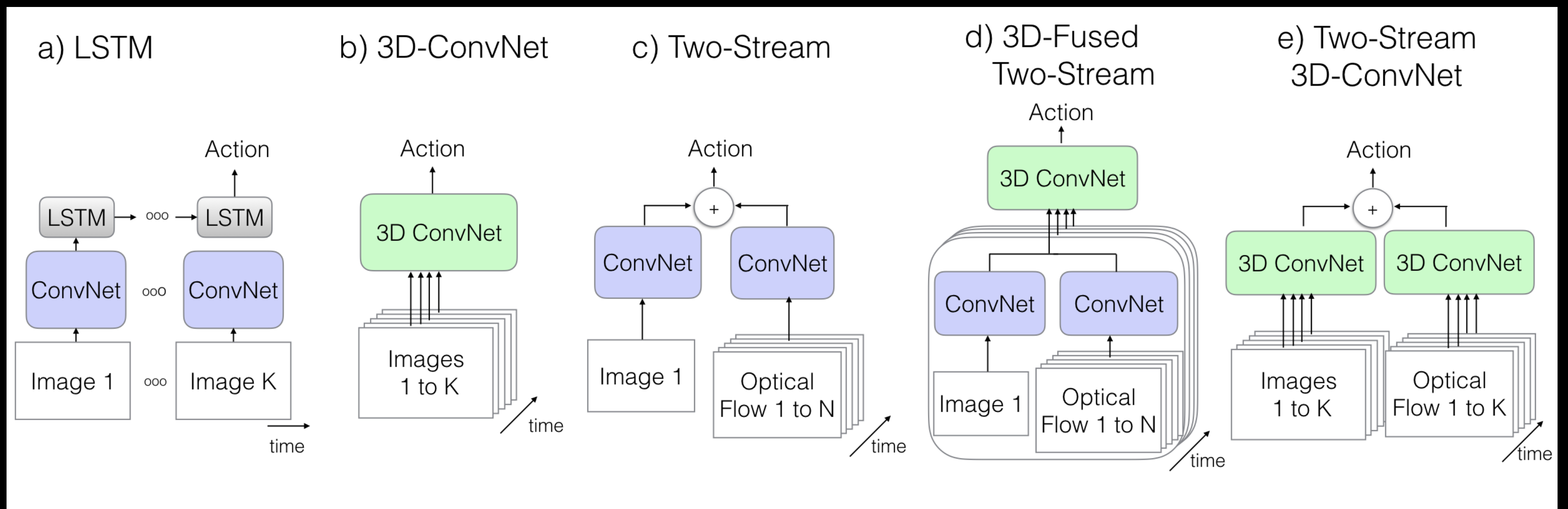
(a) Spatial ConvNet.

Training setting	Dropout ratio	
	0.5	0.9
From scratch	42.5%	52.3%
Pre-trained + fine-tuning	70.8%	72.8%
Pre-trained + last layer	72.7%	59.9%

(b) Temporal ConvNet.

Input configuration	Mean subtraction	
	off	on
Single-frame optical flow ($L = 1$)	-	73.9%
Optical flow stacking (1) ($L = 5$)	-	80.4%
Optical flow stacking (1) ($L = 10$)	79.9%	81.0%
Trajectory stacking (2) ($L = 10$)	79.6%	80.2%
Optical flow stacking (1) ($L = 10$), bi-dir.	-	81.2%

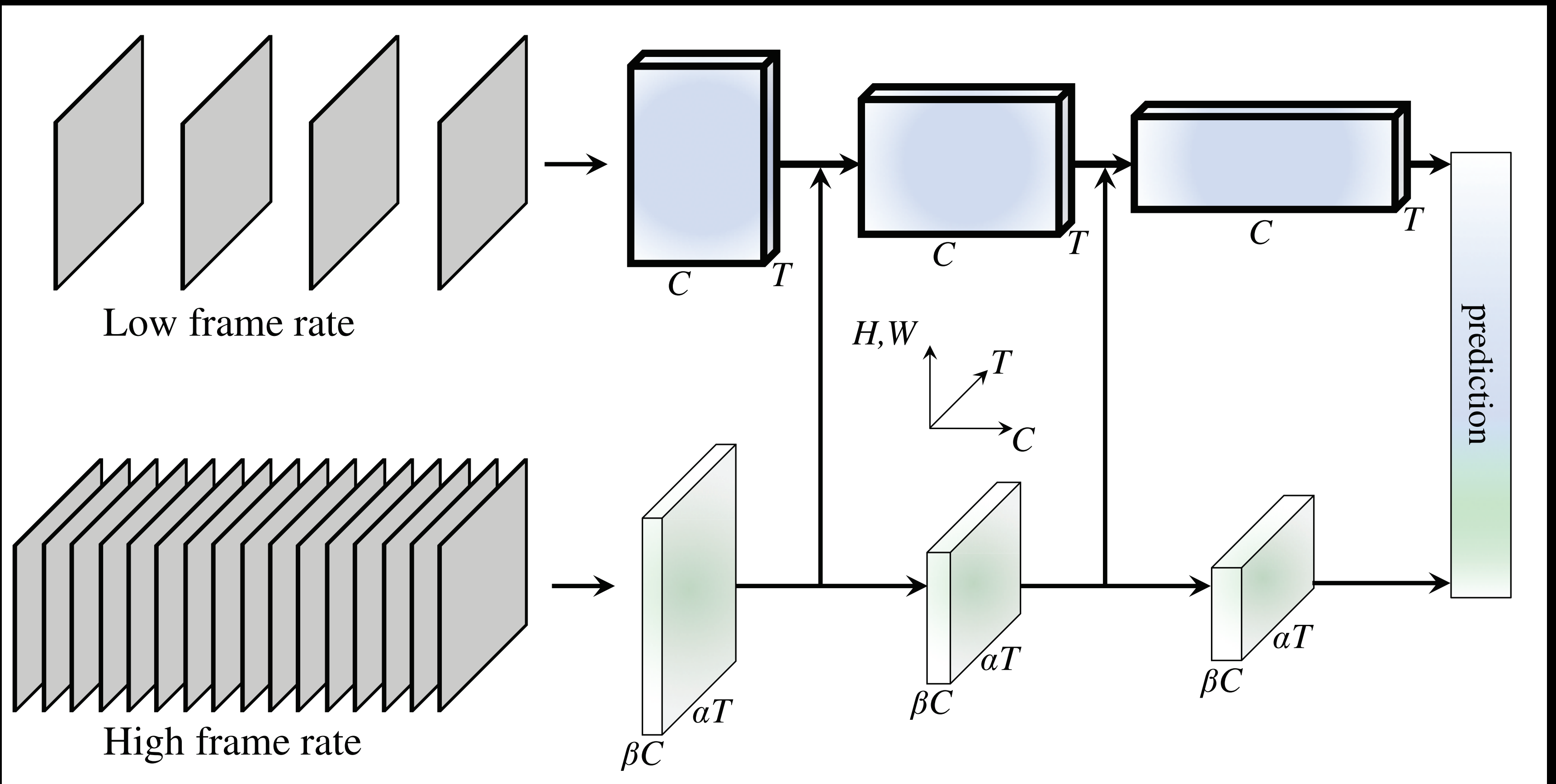
Inflated 3D Convolutions



Inflated 3D Convolutions

Architecture	UCF-101			HMDB-51			Kinetics		
	RGB	Flow	RGB + Flow	RGB	Flow	RGB + Flow	RGB	Flow	RGB + Flow
(a) LSTM	81.0	–	–	36.0	–	–	63.3	–	–
(b) 3D-ConvNet	51.6	–	–	24.3	–	–	56.1	–	–
(c) Two-Stream	83.6	85.6	91.2	43.2	56.3	58.3	62.2	52.4	65.6
(d) 3D-Fused	83.2	85.8	89.3	49.2	55.5	56.8	–	–	67.2
(e) Two-Stream I3D	84.5	90.6	93.4	49.8	61.9	66.4	71.1	63.4	74.2

SlowFast Networks



SlowFast Networks

stage	Slow pathway	Fast pathway	output sizes $T \times S^2$
raw clip	-	-	64×224^2
data layer	stride 16, 1^2	stride 2, 1^2	Slow : 4×224^2 Fast : 32×224^2
conv ₁	1×7^2 , 64 stride 1, 2^2	<u>5×7^2</u> , 8 stride 1, 2^2	Slow : 4×112^2 Fast : 32×112^2
pool ₁	1×3^2 max stride 1, 2^2	1×3^2 max stride 1, 2^2	Slow : 4×56^2 Fast : 32×56^2
res ₂	$\begin{bmatrix} 1 \times 1^2, 64 \\ 1 \times 3^2, 64 \\ 1 \times 1^2, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} \underline{3 \times 1^2}, 8 \\ 1 \times 3^2, 8 \\ 1 \times 1^2, 32 \end{bmatrix} \times 3$	Slow : 4×56^2 Fast : 32×56^2
res ₃	$\begin{bmatrix} 1 \times 1^2, 128 \\ 1 \times 3^2, 128 \\ 1 \times 1^2, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} \underline{3 \times 1^2}, 16 \\ 1 \times 3^2, 16 \\ 1 \times 1^2, 64 \end{bmatrix} \times 4$	Slow : 4×28^2 Fast : 32×28^2
res ₄	$\begin{bmatrix} 3 \times 1^2, 256 \\ 1 \times 3^2, 256 \\ 1 \times 1^2, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} \underline{3 \times 1^2}, 32 \\ 1 \times 3^2, 32 \\ 1 \times 1^2, 128 \end{bmatrix} \times 6$	Slow : 4×14^2 Fast : 32×14^2
res ₅	$\begin{bmatrix} 3 \times 1^2, 512 \\ 1 \times 3^2, 512 \\ 1 \times 1^2, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} \underline{3 \times 1^2}, 64 \\ 1 \times 3^2, 64 \\ 1 \times 1^2, 256 \end{bmatrix} \times 3$	Slow : 4×7^2 Fast : 32×7^2
global average pool, concat, fc			# classes

Table 1. An example instantiation of the SlowFast network. The dimensions of kernels are denoted by $\{T \times S^2, C\}$ for temporal, spatial, and channel sizes. Strides are denoted as $\{\text{temporal stride, spatial stride}^2\}$. Here the speed ratio is $\alpha = 8$ and the channel ratio is $\beta = 1/8$. τ is 16. The green colors mark higher temporal resolution, and orange colors mark fewer channels, for the Fast pathway. Non-degenerate temporal filters are underlined. Residual blocks are shown by brackets. The backbone is ResNet-50.

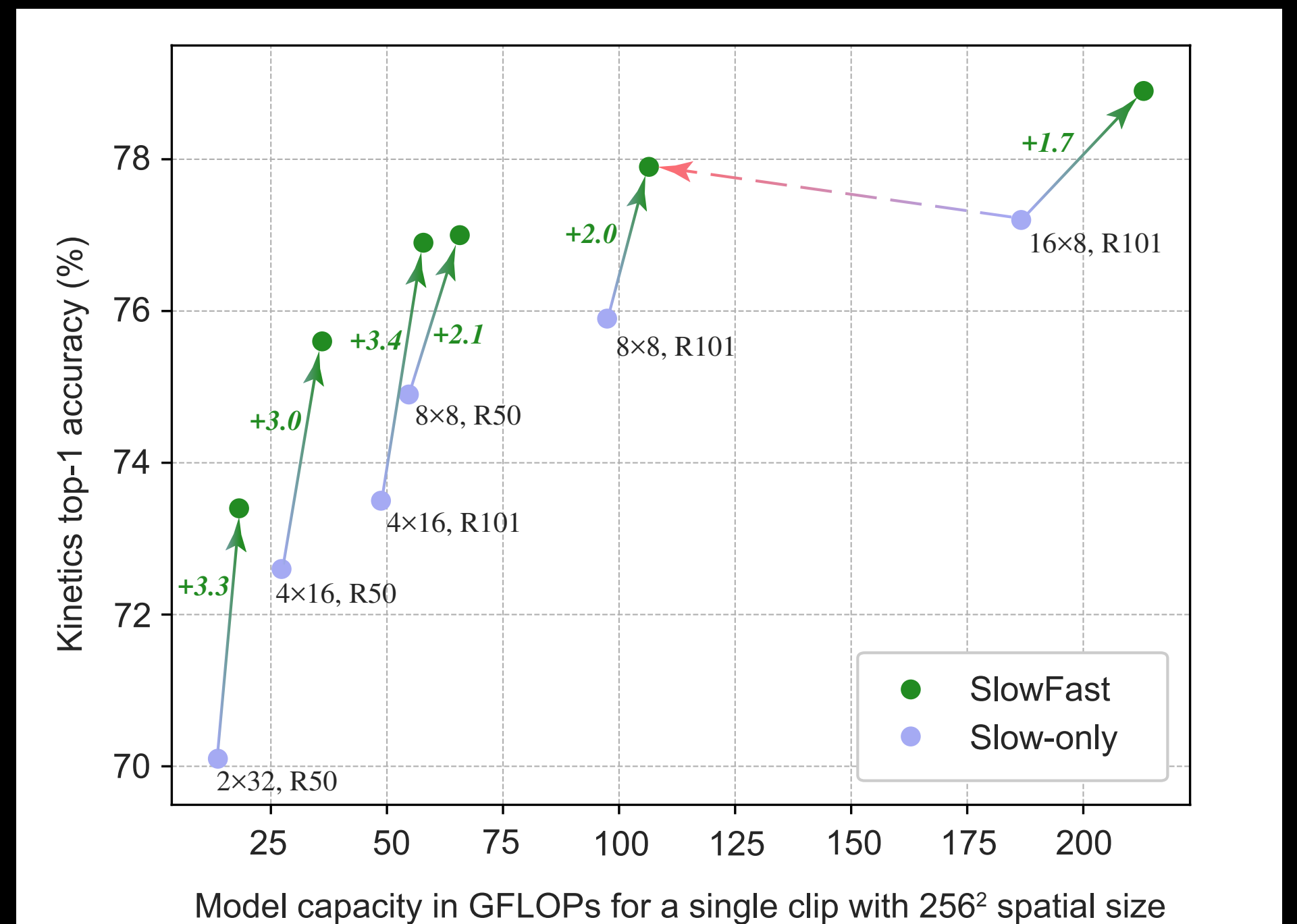
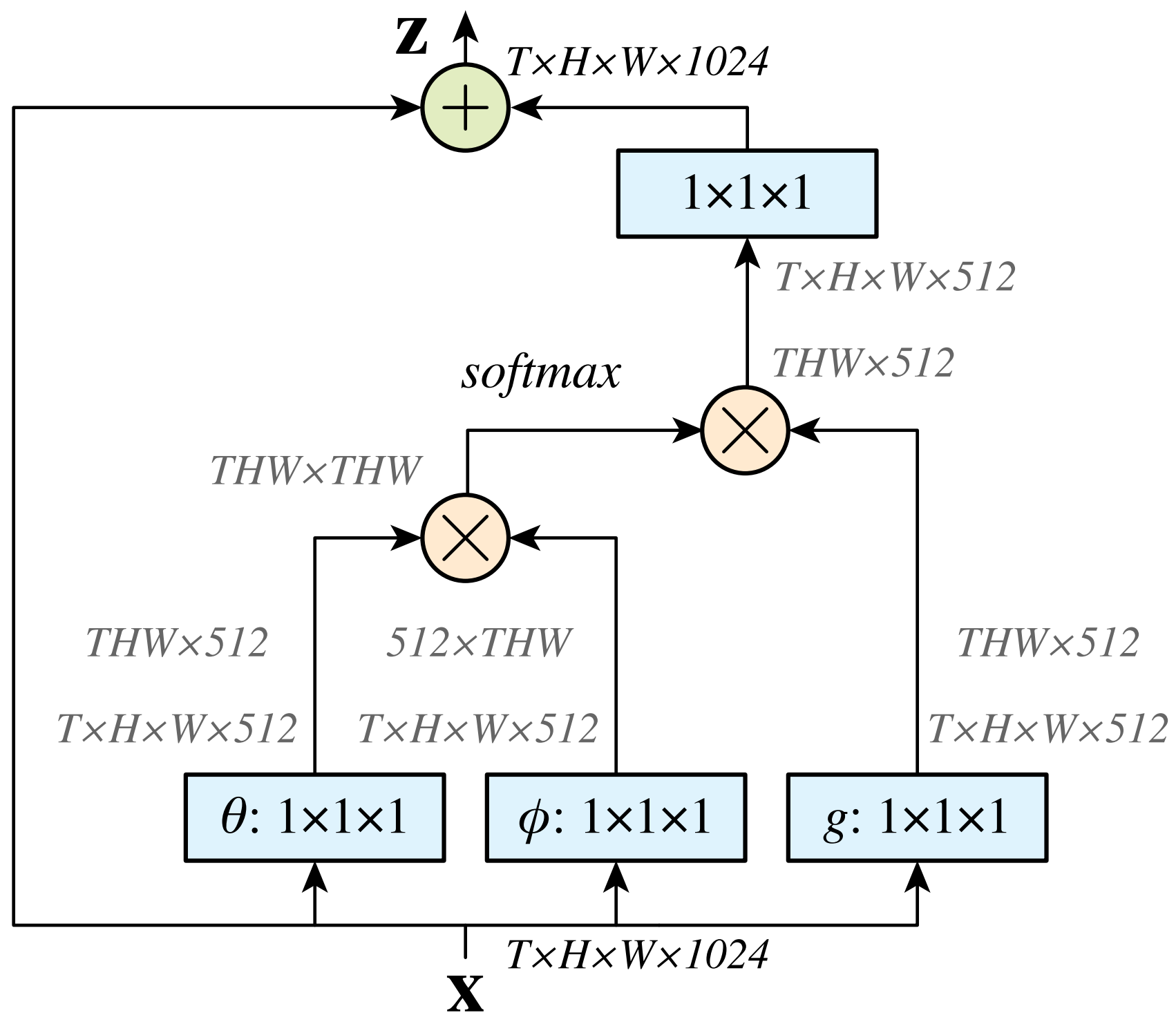


Figure 2. Accuracy/complexity tradeoff on Kinetics-400 for the SlowFast (green) vs. Slow-only (blue) architectures. SlowFast is consistently better than its Slow-only counterpart in all cases (green arrows). SlowFast provides higher accuracy and lower cost than temporally heavy Slow-only (e.g. red arrow). The complexity is for a single 256^2 view, and accuracy are obtained by 30-view testing.

Non-local Networks



$$\mathbf{y}_i = \frac{1}{\mathcal{C}(\mathbf{x})} \sum_{\forall j} f(\mathbf{x}_i, \mathbf{x}_j) g(\mathbf{x}_j).$$

Non-local Networks



Outline

- Correspondence Problems
 - Optical Flow
 - Tracking
 - Mid-level Correspondence
- Recognition in Videos
- Videos as a source of supervision

Videos as a source for supervision

