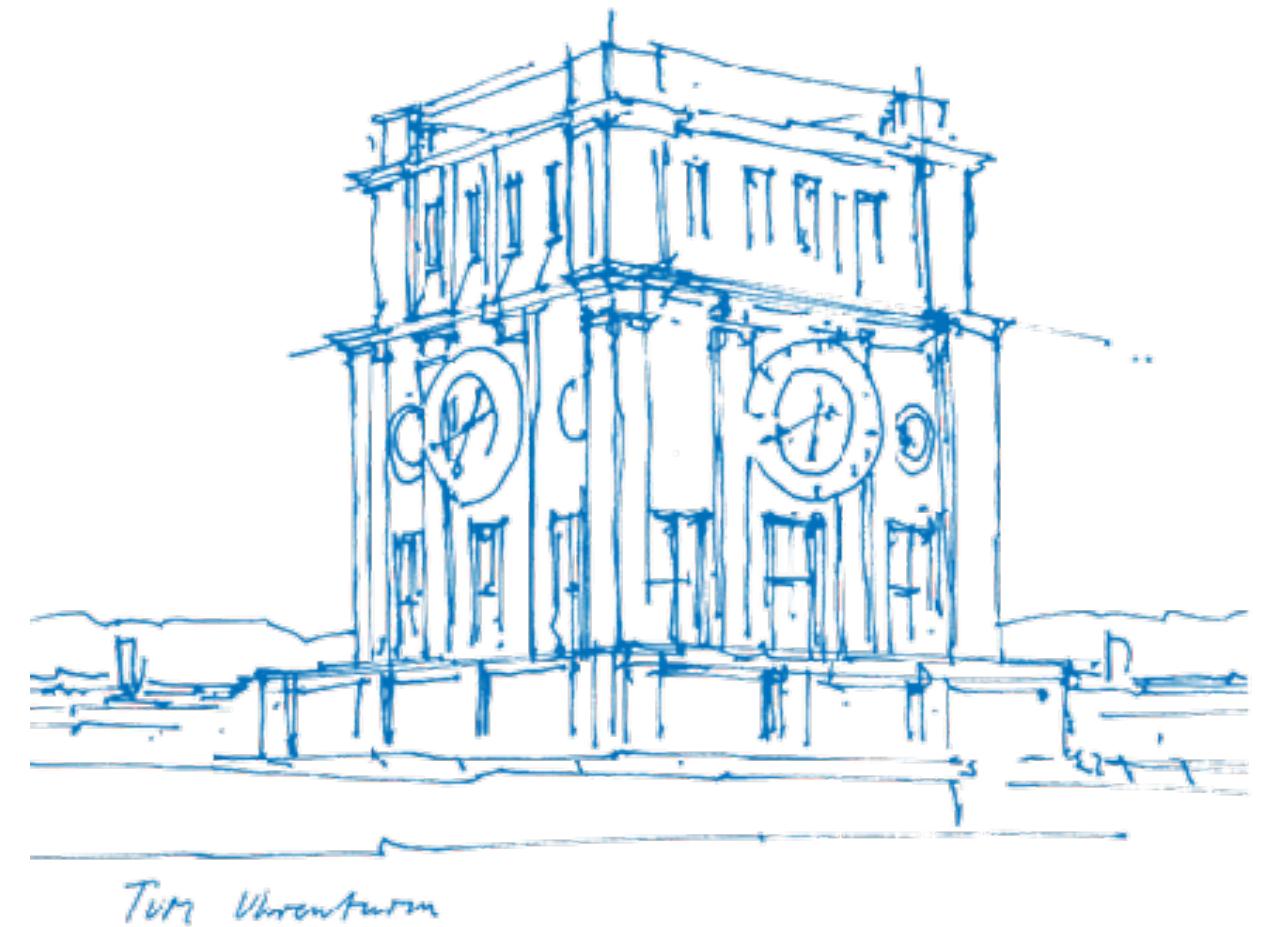


# Computer Vision III:

## MOT 02: Offline tracking

Nikita Araslanov  
29.11.2022

Content credit:  
Prof. Laura Leal-Taixé  
<https://dvl.in.tum.de>



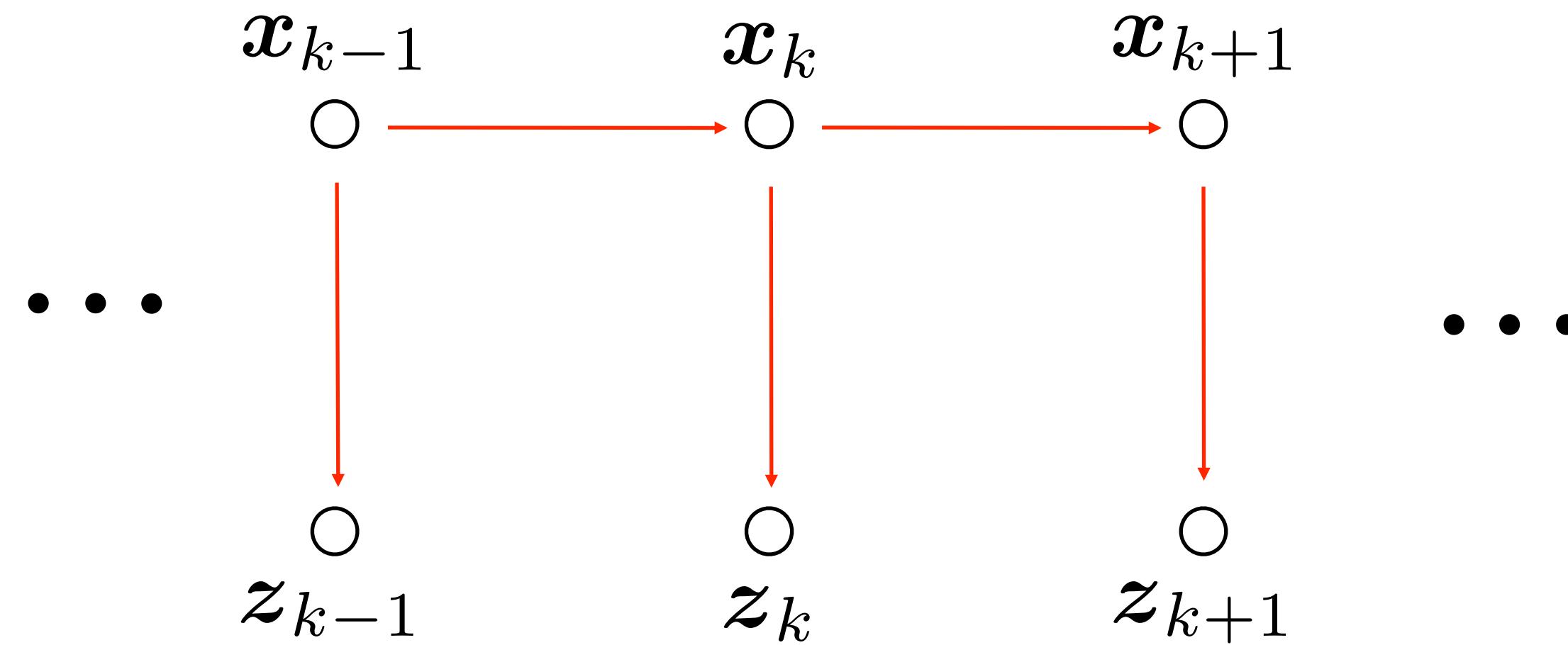
# Problem statement

- Our focus: object detection and temporal association



# Bayesian graphical model

- Hidden Markov model:



Assumptions:

$$p(z_k | x_k, \mathbf{Z}_{k-1}) = p(z_k | x_k) \quad p(x_k | x_{k-1}, \mathbf{Z}_{k-1}) = p(x_k | x_{k-1})$$

$$p(x_k | \mathbf{X}_{k-1}) = p(x_k | x_{k-1})$$

[Stefan Roth]

# Bayesian formulation

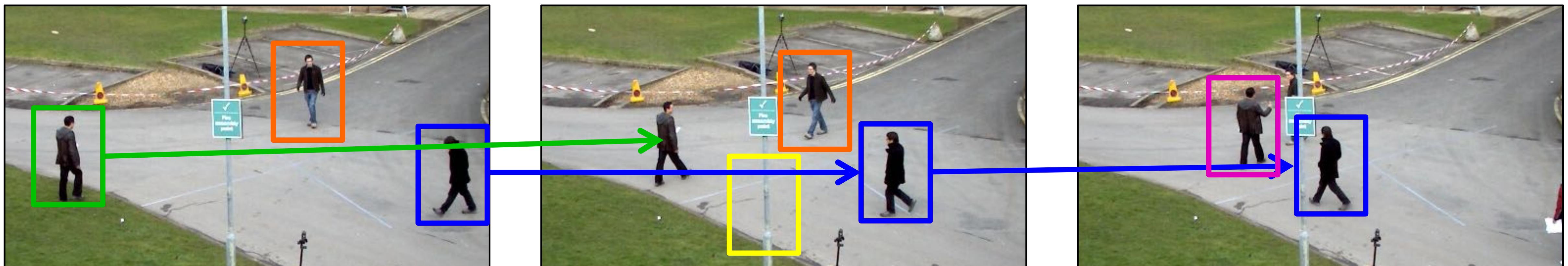
$$p(\mathbf{x}_k | \mathbf{Z}_k) = \kappa \cdot p(z_k | \mathbf{x}_k) \cdot \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) \cdot p(\mathbf{x}_{k-1} | \mathbf{Z}_{k-1}) d\mathbf{x}_{k-1}$$

$p(\mathbf{x}_k   \mathbf{Z}_k)$	posterior probability at current time step
$p(z_k   \mathbf{x}_k)$	likelihood $\rightarrow$ How will current estimate of tracker
$p(\mathbf{x}_k   \mathbf{x}_{k-1})$	<u>temporal prior</u>
$p(\mathbf{x}_{k-1}   \mathbf{Z}_{k-1})$	<u>posterior probability at previous time step</u>
$\kappa$	normalizing term

[Stefan Roth]

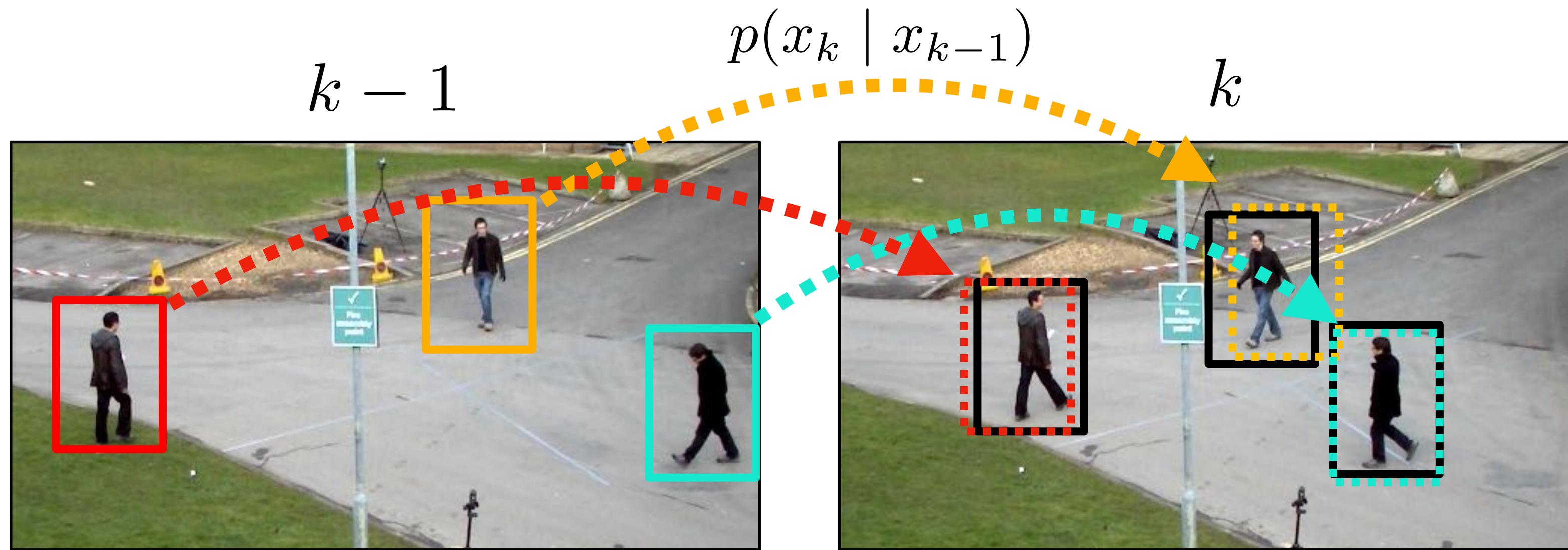
# Tracking-by-detection

- We will focus on algorithms where a set of detections is provided.
- Remember detections are not perfect!



Find detections that match and form a trajectory

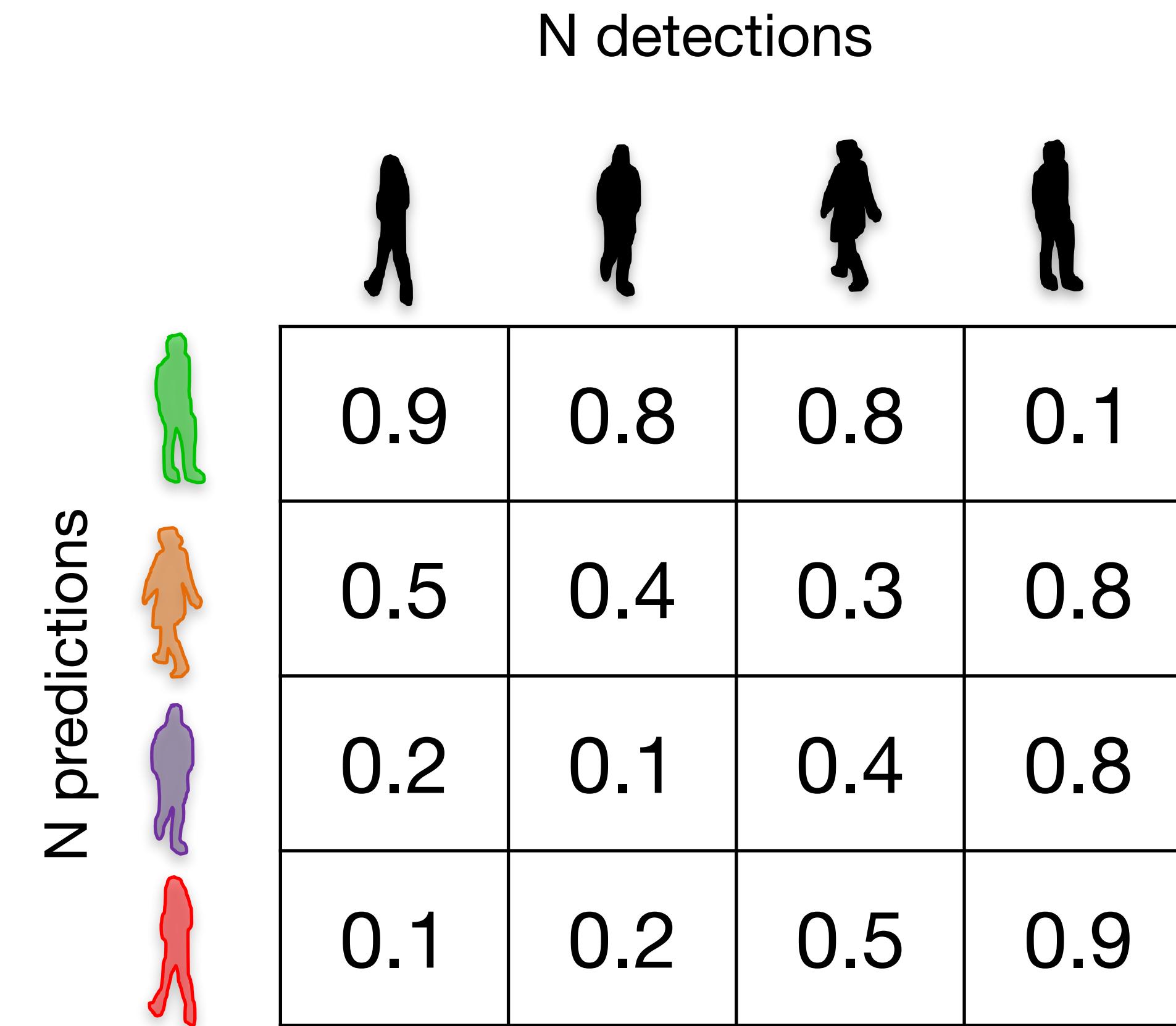
# Approach 1



- 1. Track initialization (e.g. using a detector)
- 2. Prediction of the next position (motion model)
- 3. **Matching** predictions with detections (appearance model)

# Bipartite matching

1. Define distances between boxes (e.g.,  $1 - IoU$ );



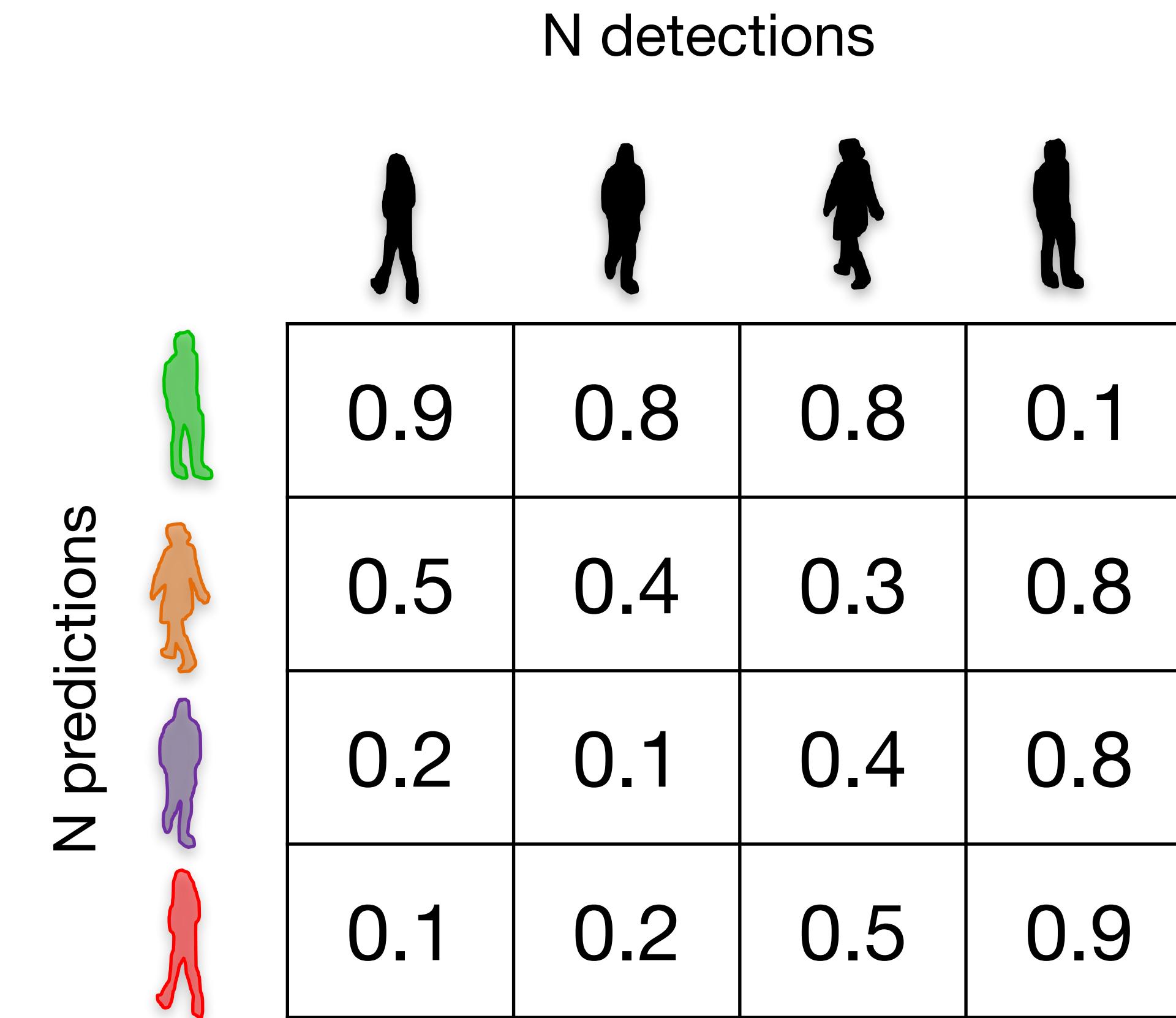
# Bipartite matching

1. Define distances between boxes (e.g.,  $1 - IoU$ );

- we obtain NxN matrix

2. Solve assignment problem

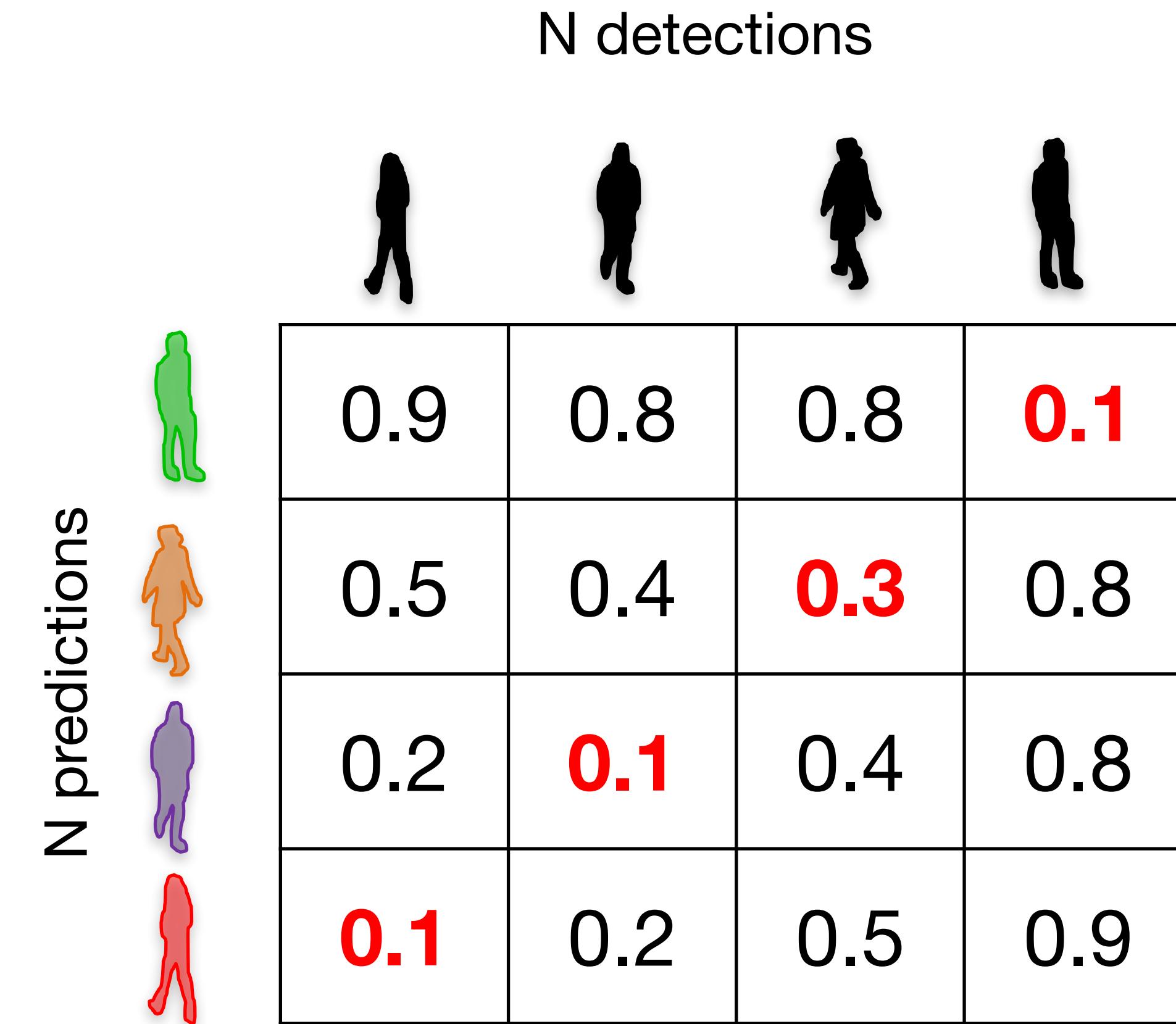
- using Hungarian algorithm\*



\*Demo: <http://www.hungarianalgorithm.com/solve.php>

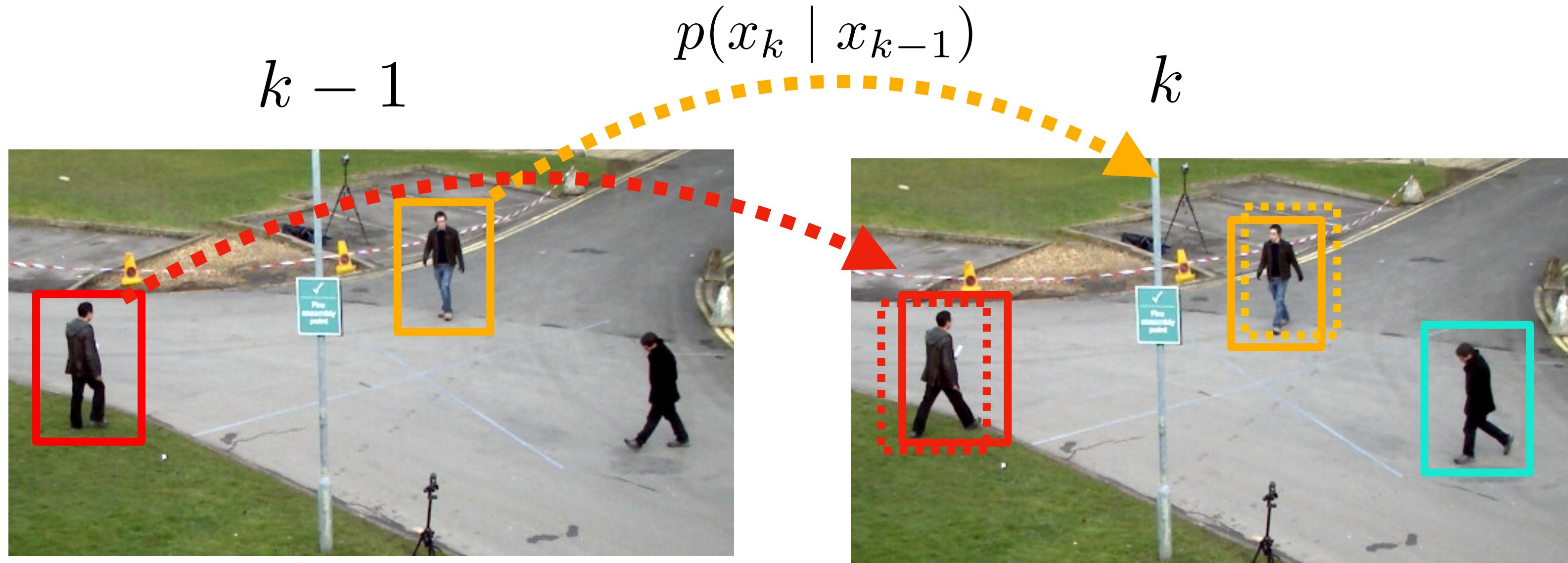
# Bipartite matching

1. Define distances between boxes (e.g.,  $1 - IoU$ ):
  - we obtain NxN matrix
2. Solve assignment problem
  - using Hungarian algorithm\*  $O(N^3)$
3. The bipartite matching solution corresponds to the minimum total cost



\*Demo: <http://www.hungarianalgorithm.com/solve.php>

# Approach 2



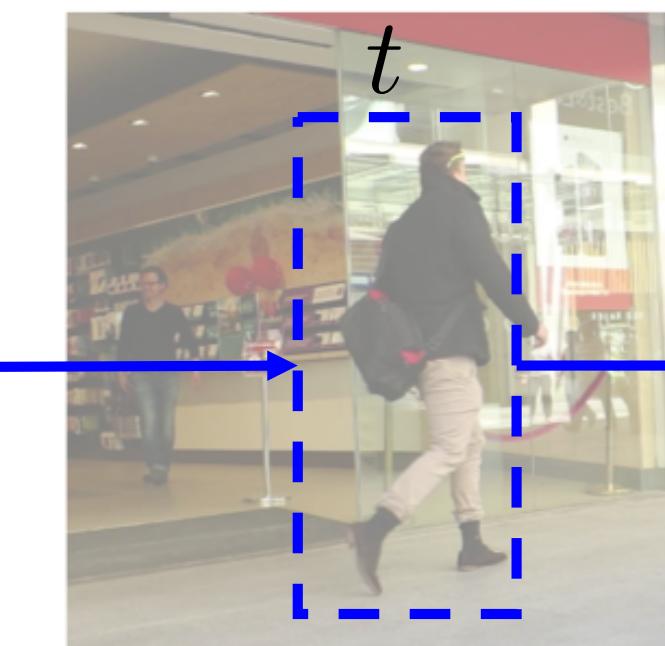
1. Detect objects in frame  $k-1$  (e.g. using an object detector)
2. Initialise in the same location in frame  $k$
3. **Refine** predictions with **regression**
4. Run object detection again to find new tracks

# Approach 2: Tracktor

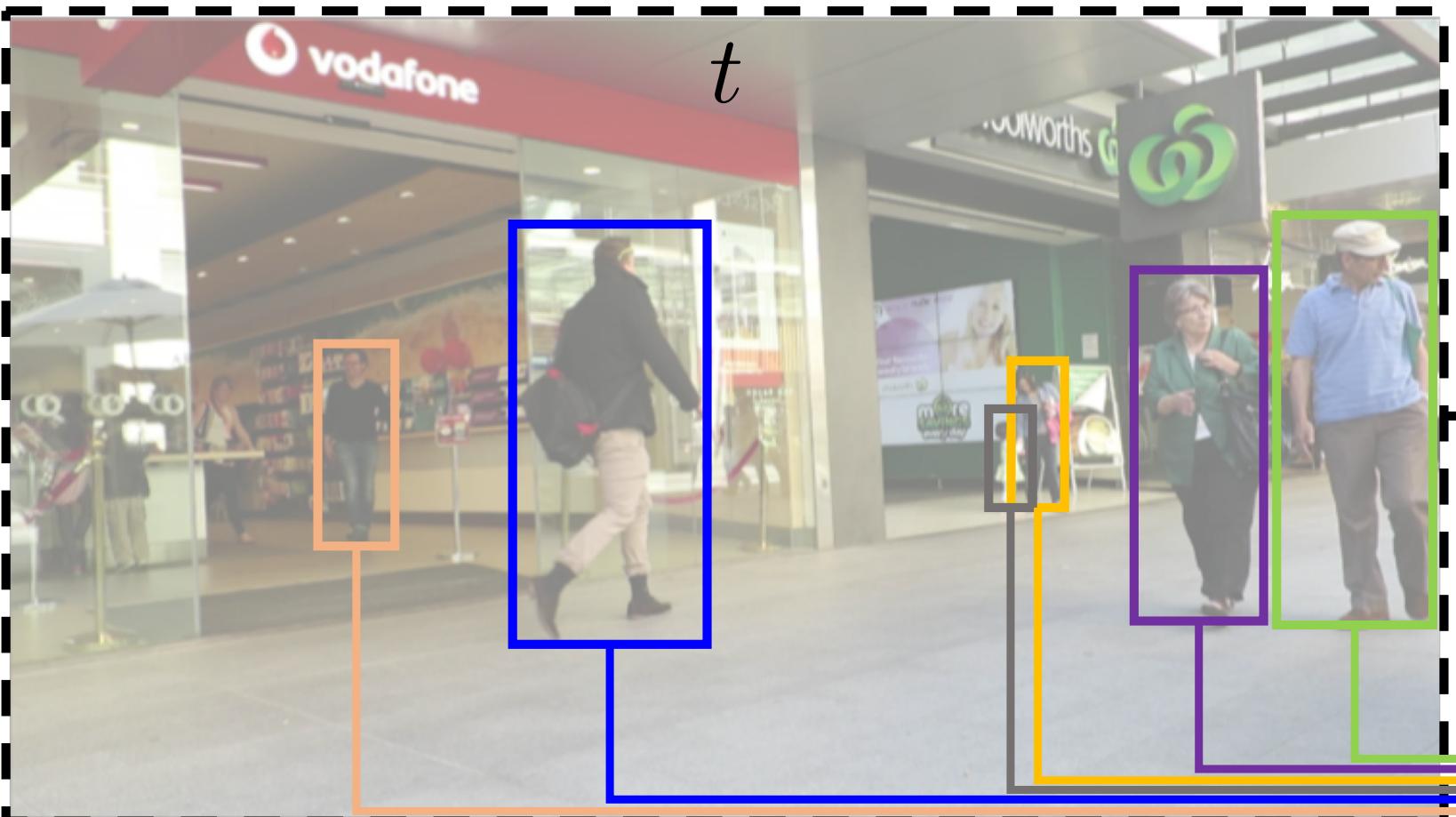
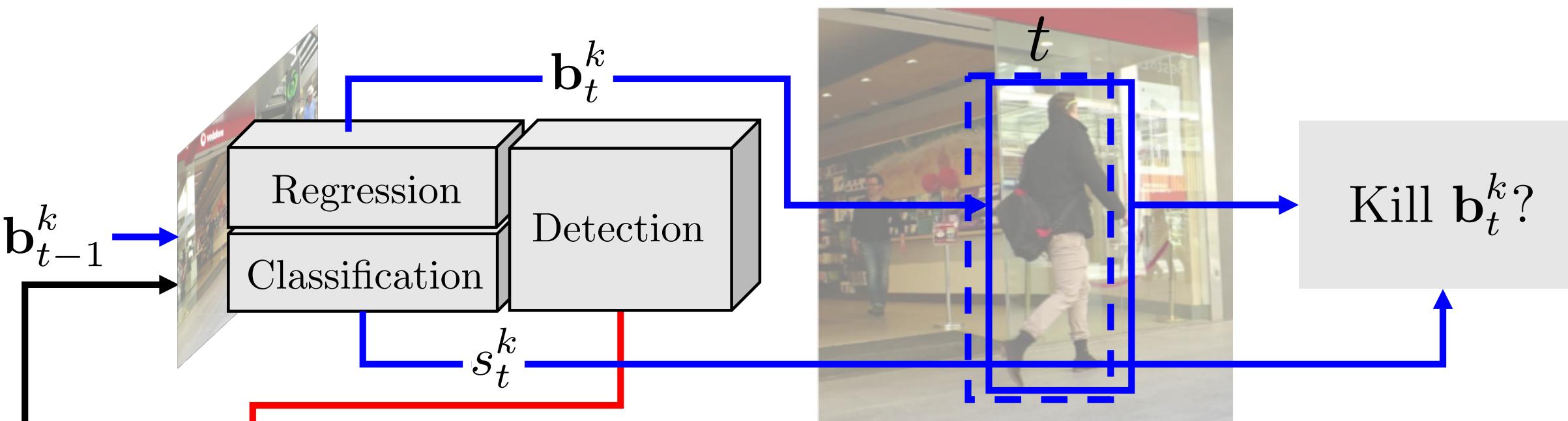
1. Run detection



2. Copy boxes the next frame



3. Refine boxes with regression

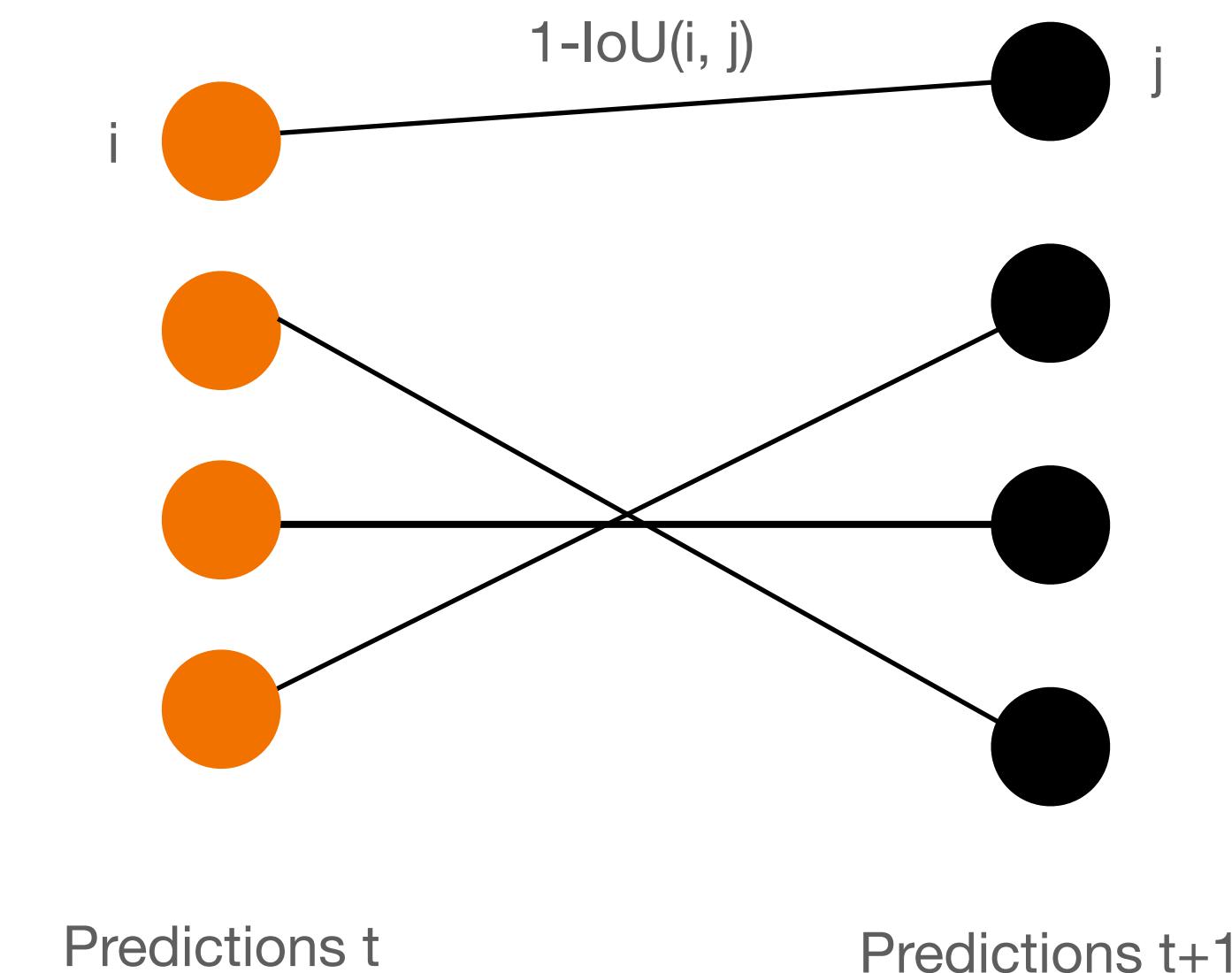


4. Initialise new tracks

Bergmann et al., "Tracking without bells and whistles". ICCV 2019

# Open questions

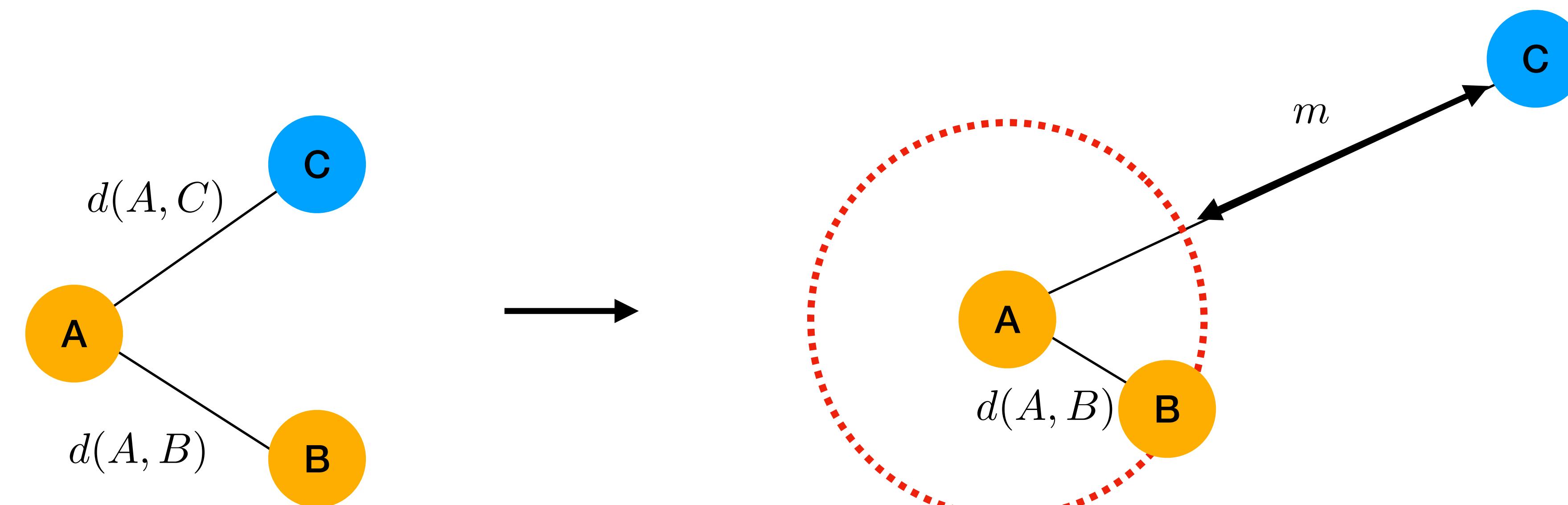
- How to match motion prediction with detections?
  - **in general: reduction to the assignment problem;**
  - small motion: refine the boxes with a regression head
- We have used IoU so far
  - implicit assumption of small motion
- We need a more robust **metric!**



# Metric learning: triplet loss

$$\mathcal{L}(A, B, C) = \max(0, \|f(A) - f(B)\|^2 - \|f(A) - f(C)\|^2 + m)$$

Intuitive idea:



# Online vs. offline tracking

Last lecture

- Online tracking
  - Processes two frames at a time
  - For real-time applications
  - Prone to drifting → hard to recover from errors or occlusions

- Offline tracking
  - Processes a batch of frames
  - Good to recover from occlusions (short ones as we will see)
  - Not suitable for real-time applications
  - Suitable for video analysis

# Online vs. offline tracking

- Online tracking
  - Processes two frames at a time
  - For real-time applications
  - Prone to drifting → hard to recover from errors or occlusions

- Offline tracking
  - Processes a batch of frames
  - Good to recover from occlusions (short ones as we will see)
  - Not suitable for real-time applications
  - Suitable for video analysis

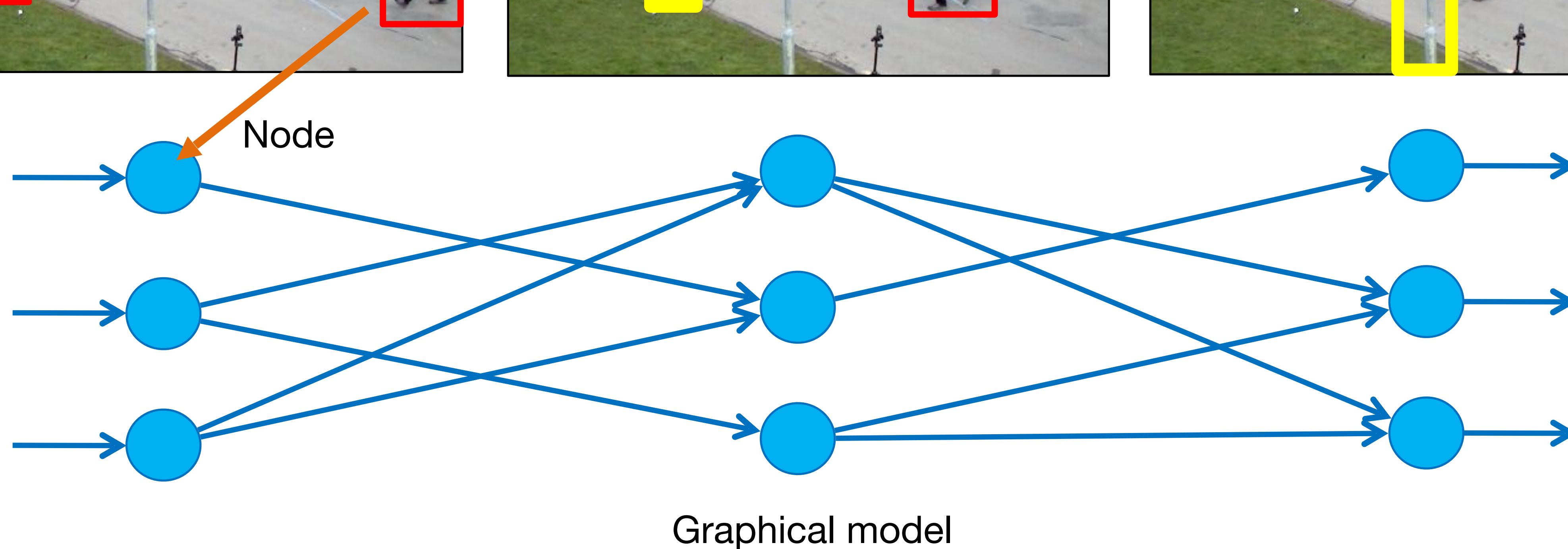
This lecture

→ Large context

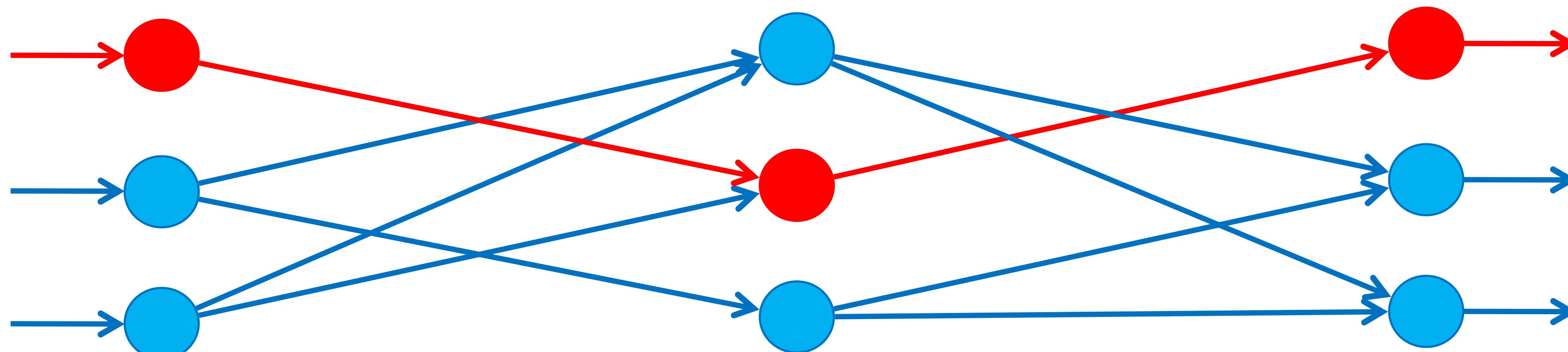
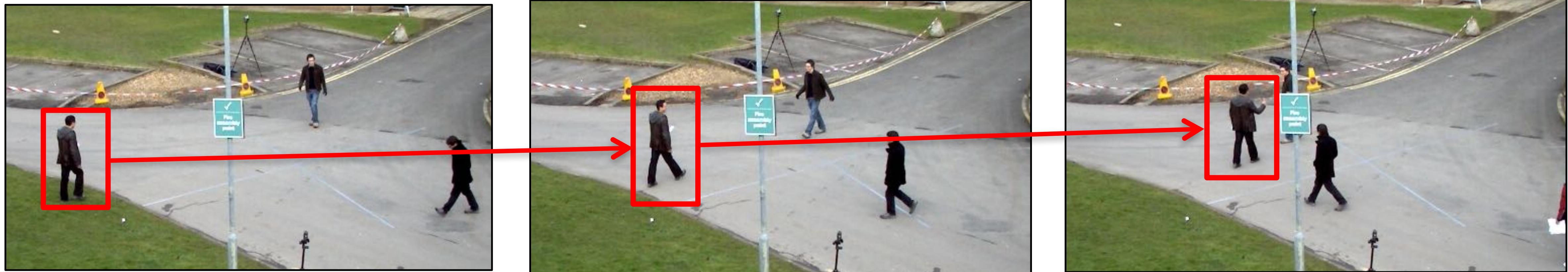
# Graph-based MOT



# Tracking with network flows

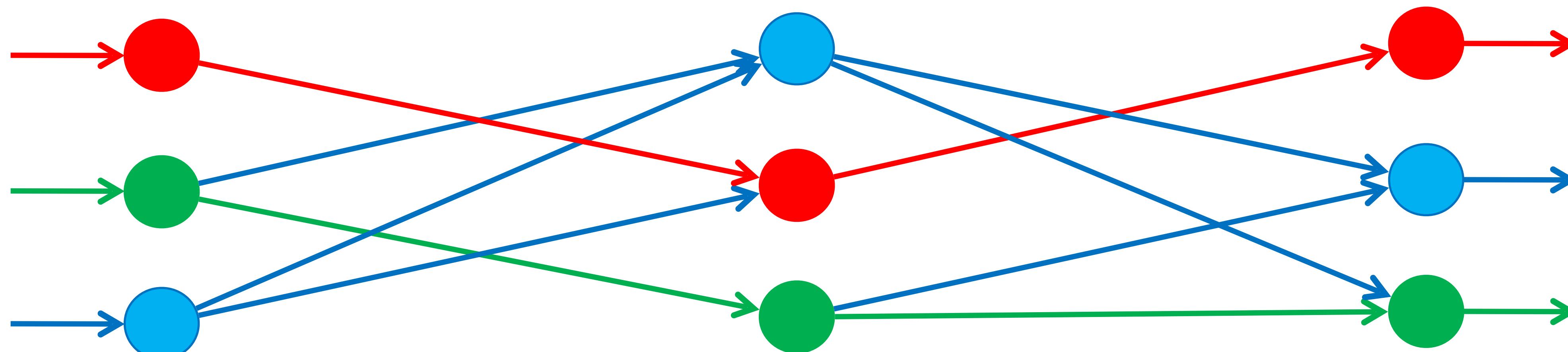
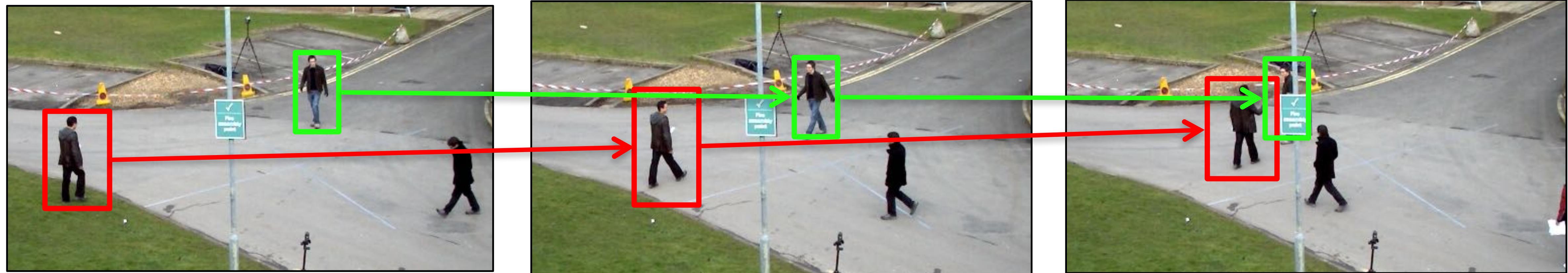


# Tracking with network flows



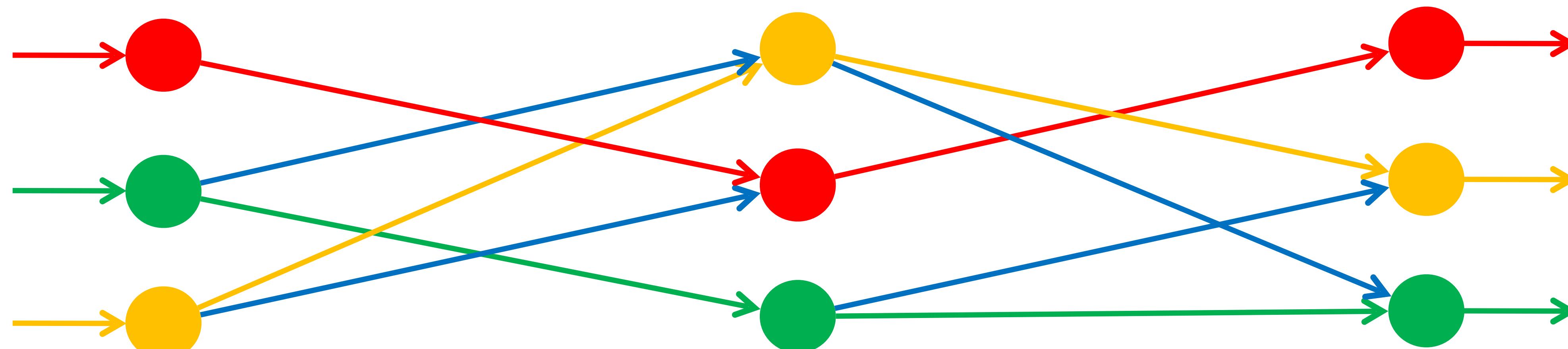
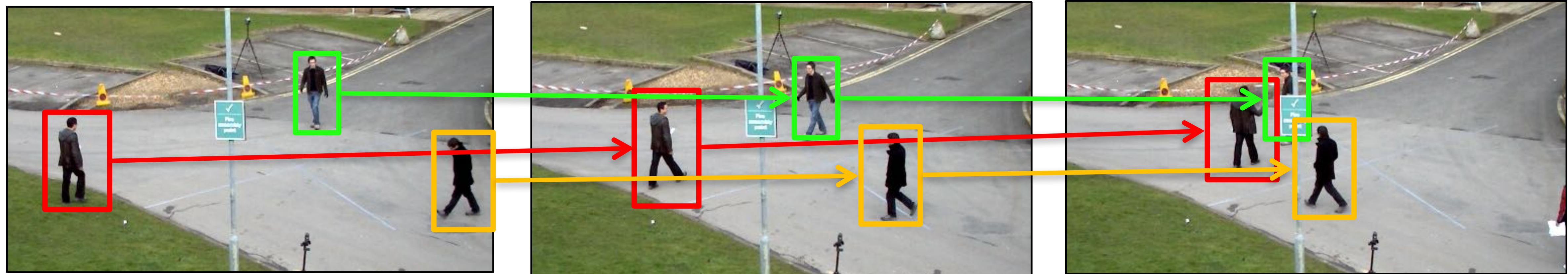
Zhang et al. "Global Data Association for Multi-Object Tracking Using Network Flows". CVPR 2008

# Tracking with network flows



Zhang et al. "Global Data Association for Multi-Object Tracking Using Network Flows". CVPR 2008

# Tracking with network flows



Zhang et al. "Global Data Association for Multi-Object Tracking Using Network Flows". CVPR 2008

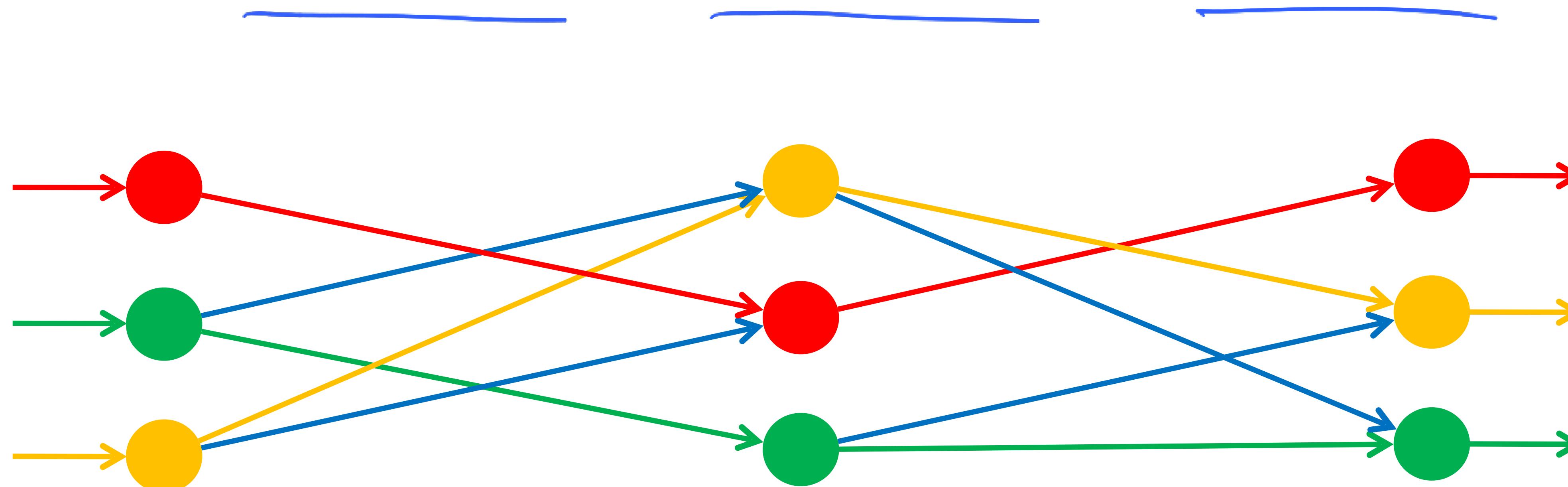
# Tracking with network flows

- Node = object detection
- Edge = flow = trajectory
- 1 unit of flow = 1 pedestrian
- Goal: disjoint set of trajectories

# Tracking with network flows

- Minimum-cost maximum-flow problem

“Determine the maximum flow with a minimum cost”

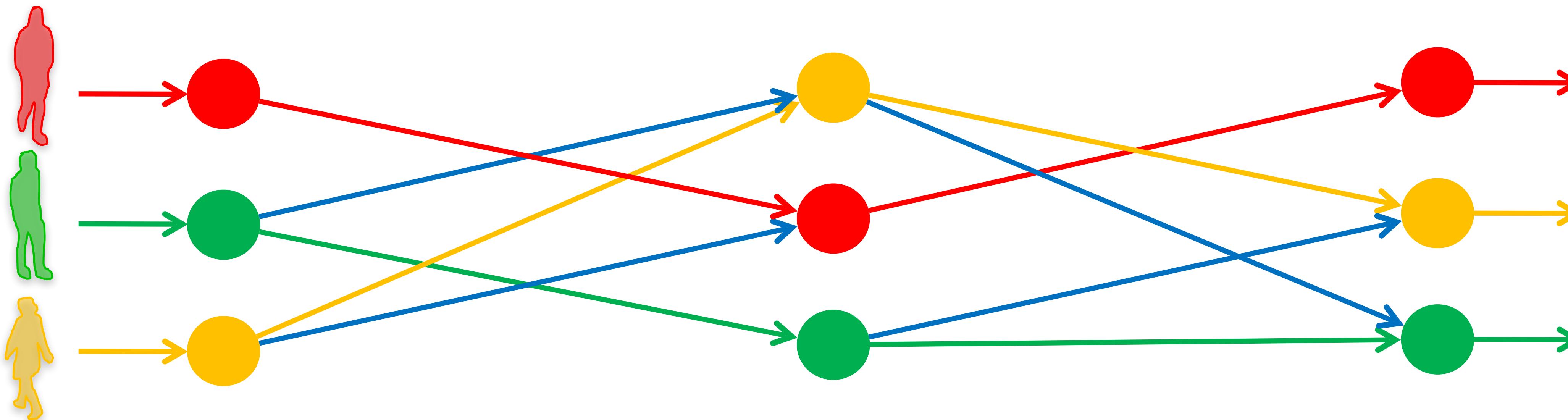


Ravindra K. Ahuja et al., “Network Flows: Theory, Algorithms, and Applications” (1993).

# Tracking with network flows

- Minimum-cost maximum-flow problem

“Determine the maximum flow with a minimum cost”

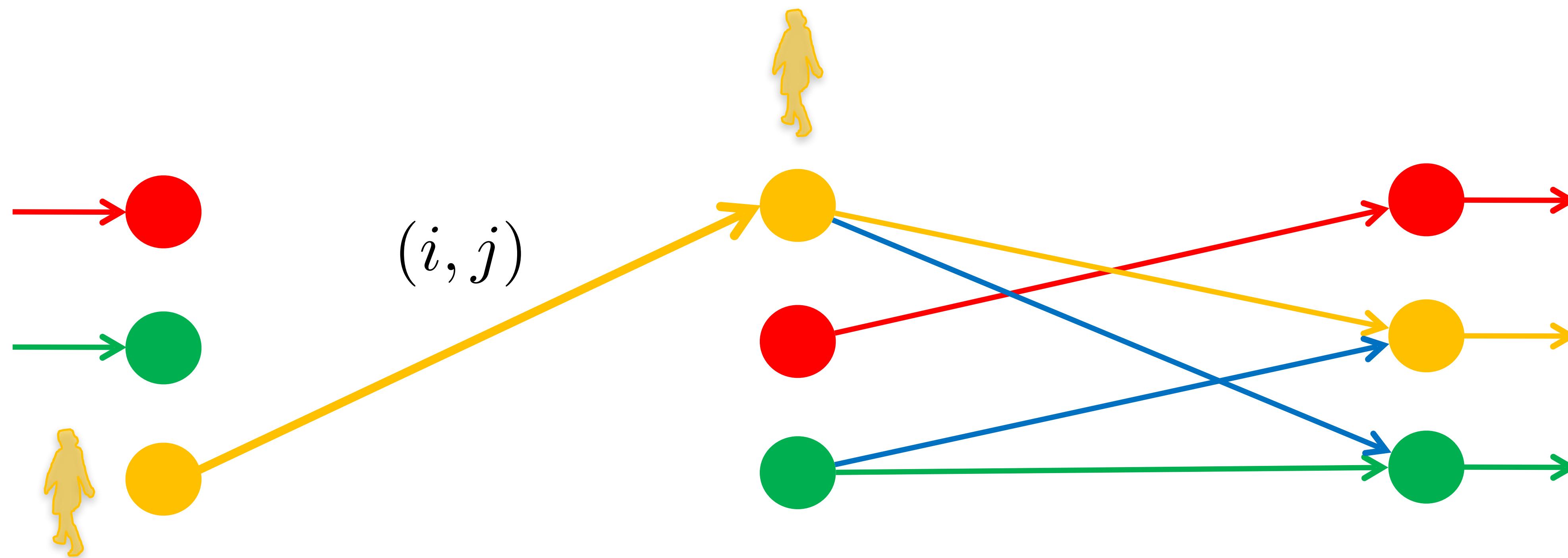


Ravindra K. Ahuja et al., “Network Flows: Theory, Algorithms, and Applications” (1993).

# Tracking with network flows

- Minimising the cost:

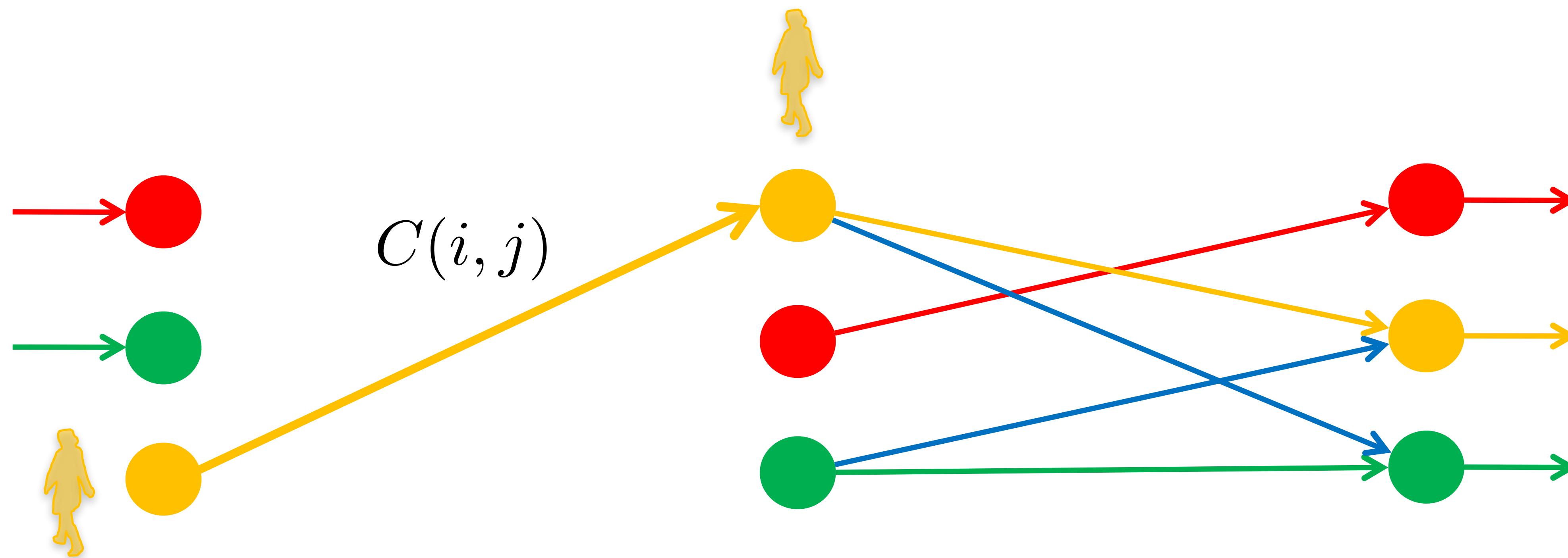
$$f^* = \arg \min_f \sum_{i,j} C(i,j) f(i,j)$$



# Tracking with network flows

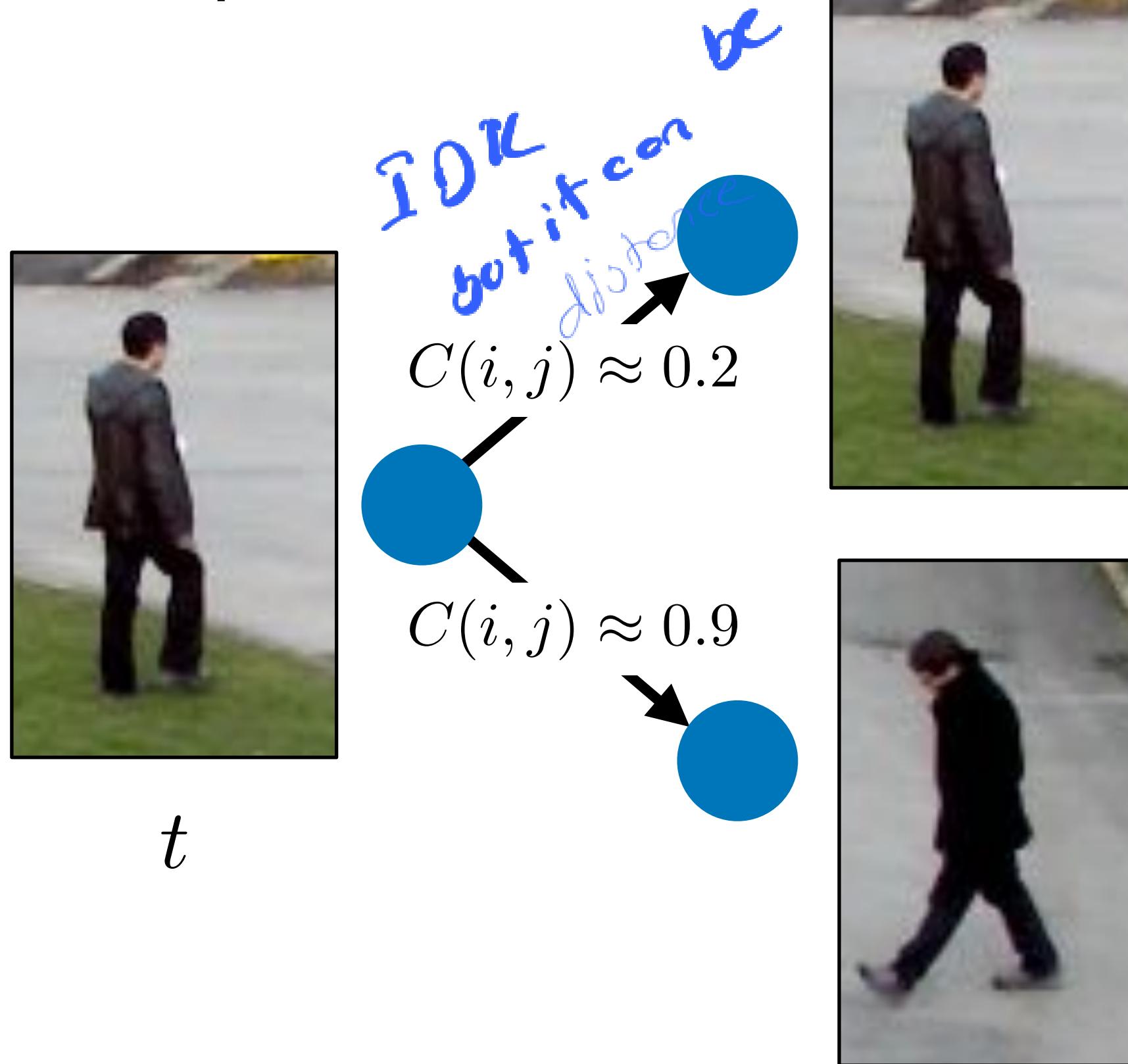
- Minimising the cost:

$$f^* = \arg \min_f \sum_{i,j} C(i,j) f(i,j)$$



# Tracking with network flows

- Minimising the cost:
- Example costs:

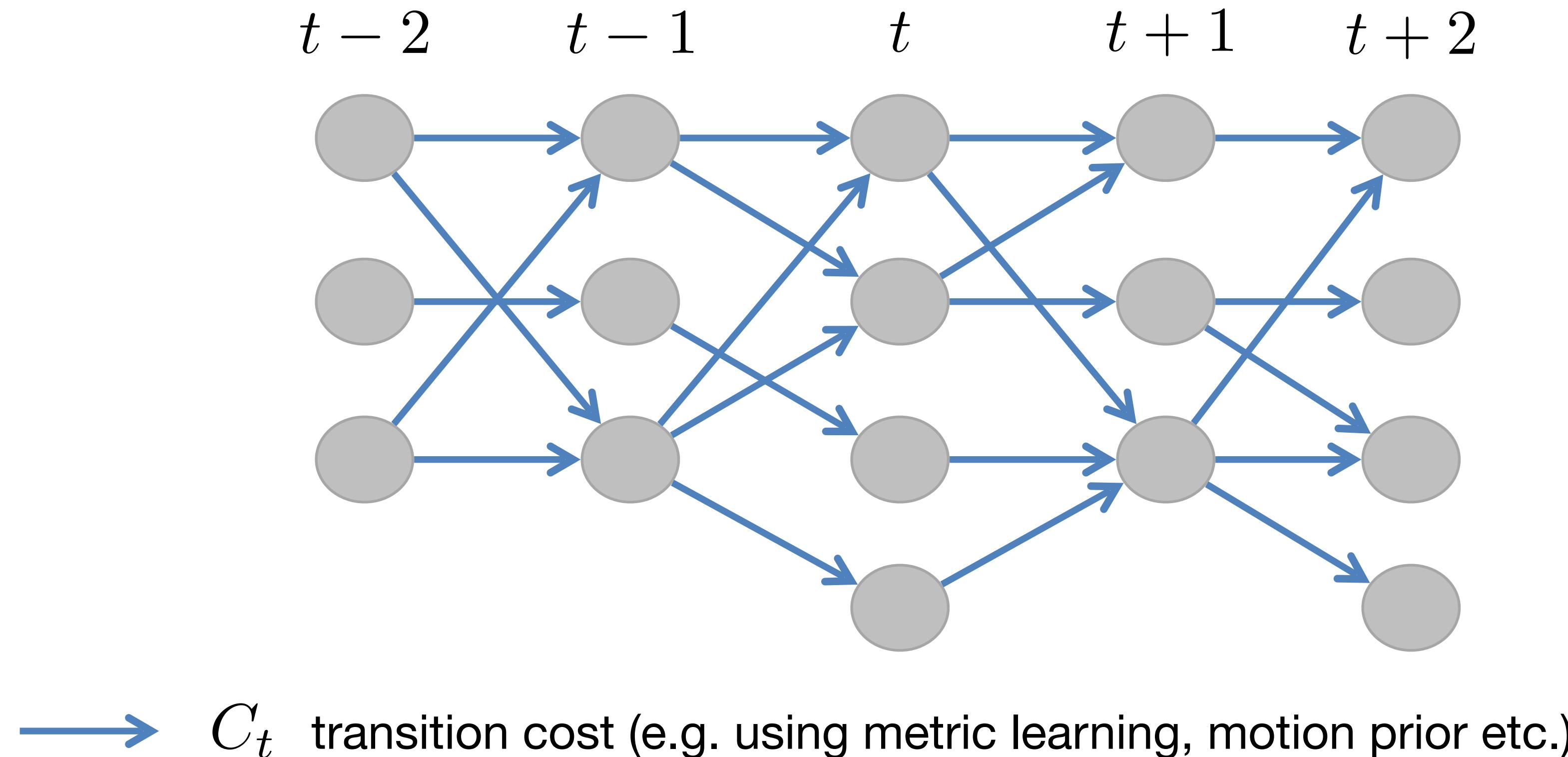


$$f^* = \arg \min_f \sum_{i,j} C(i,j) f(i,j)$$

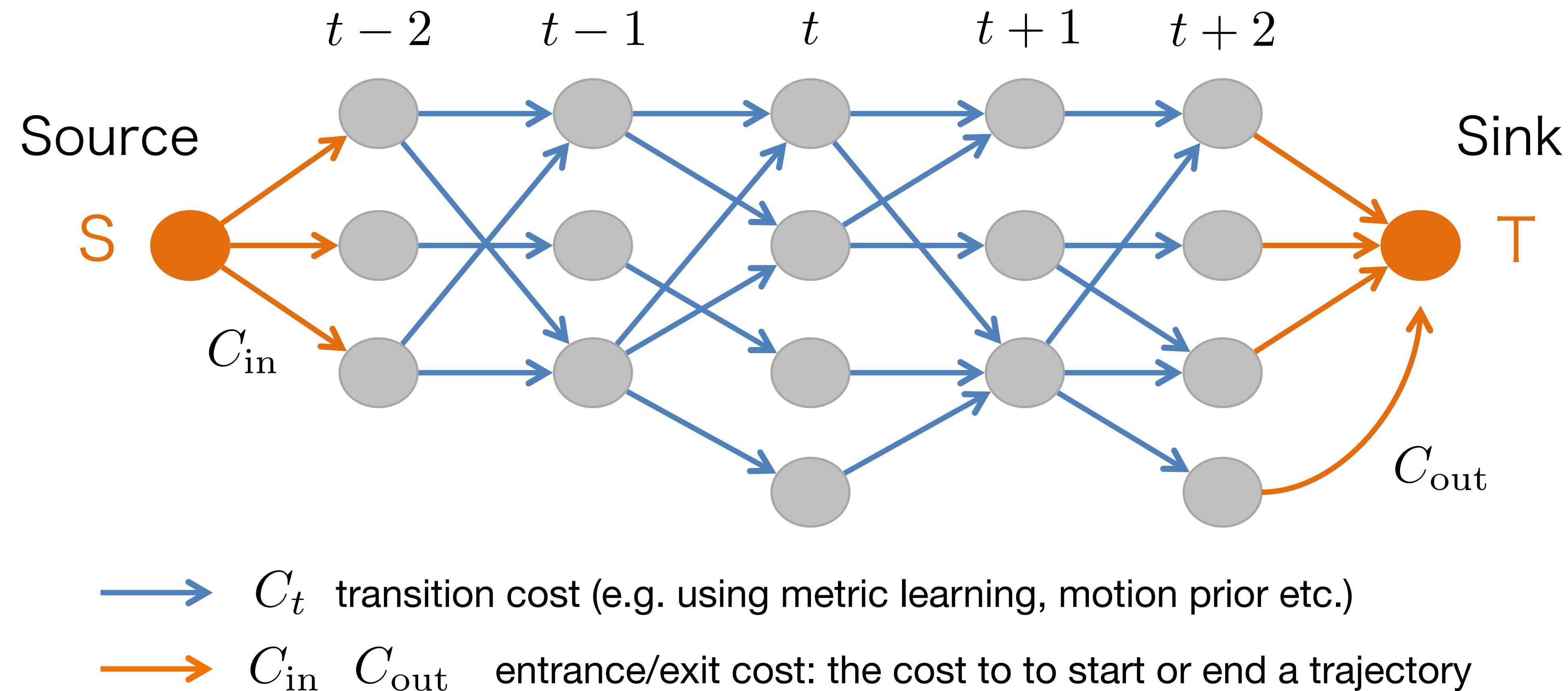
↑  
Disjoint set of  
trajectories

Cost  
 of association  
 of notes  
 Indicator  $\{0,1\}$

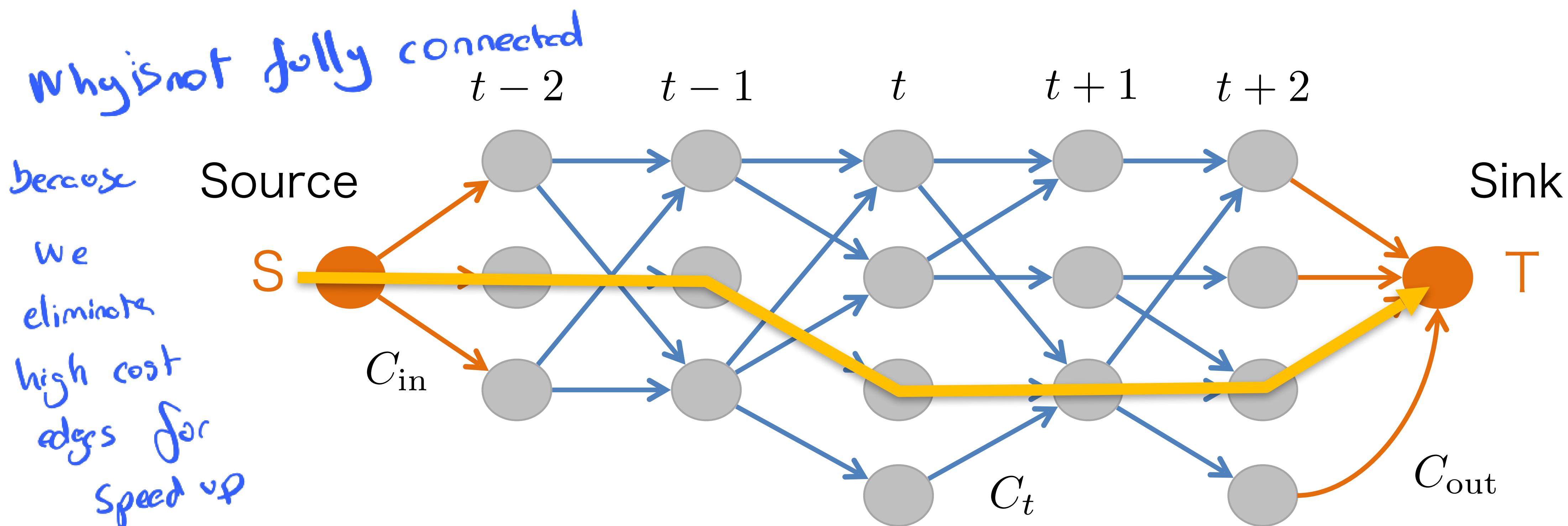
# Constructing cost-flow network



# Constructing cost-flow network



# Constructing cost-flow network

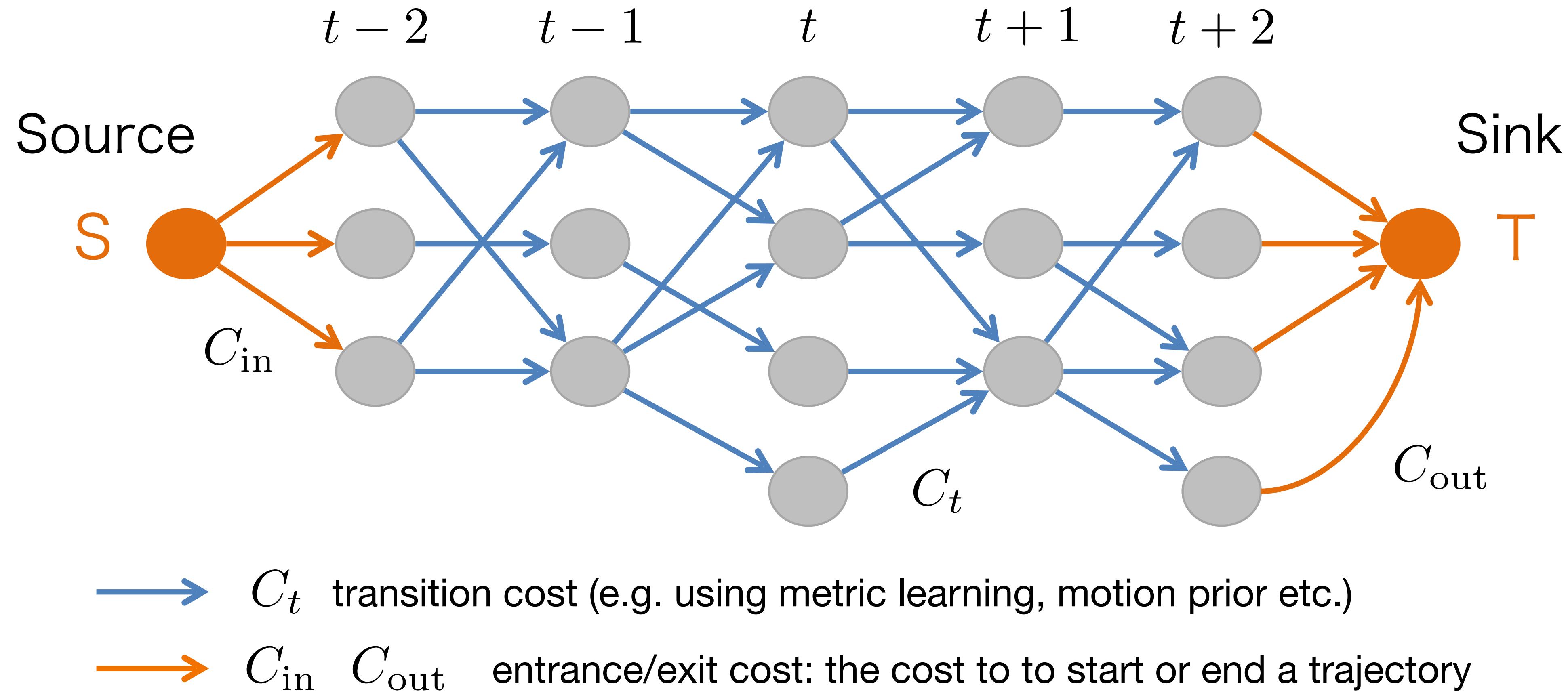


→  $C_t$  transition cost (e.g. using metric learning, motion prior etc.)

→  $C_{\text{in}} \quad C_{\text{out}}$  entrance/exit cost: the cost to start or end a trajectory

Trajectory: a path starting at  $S$  and ending at  $T$ .

# Constructing cost-flow network



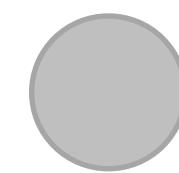
Open question: how to incorporate detection confidence?

# Tracking with network flows

- Split the node in two.
- Assign the cost using detection confidence.



before:



after:

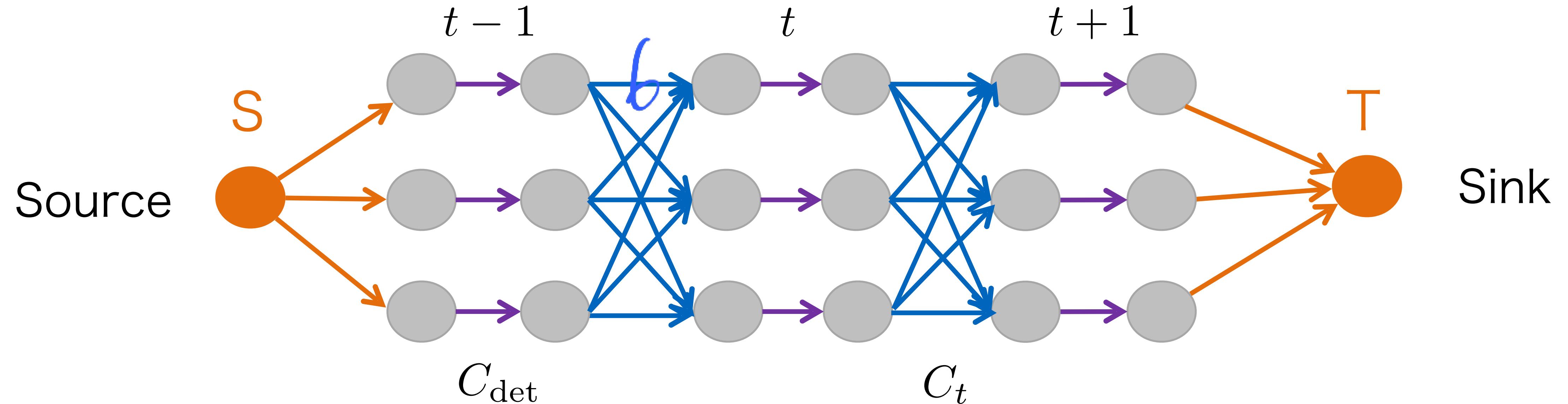


$$C_{\text{det}}$$

“detection” cost

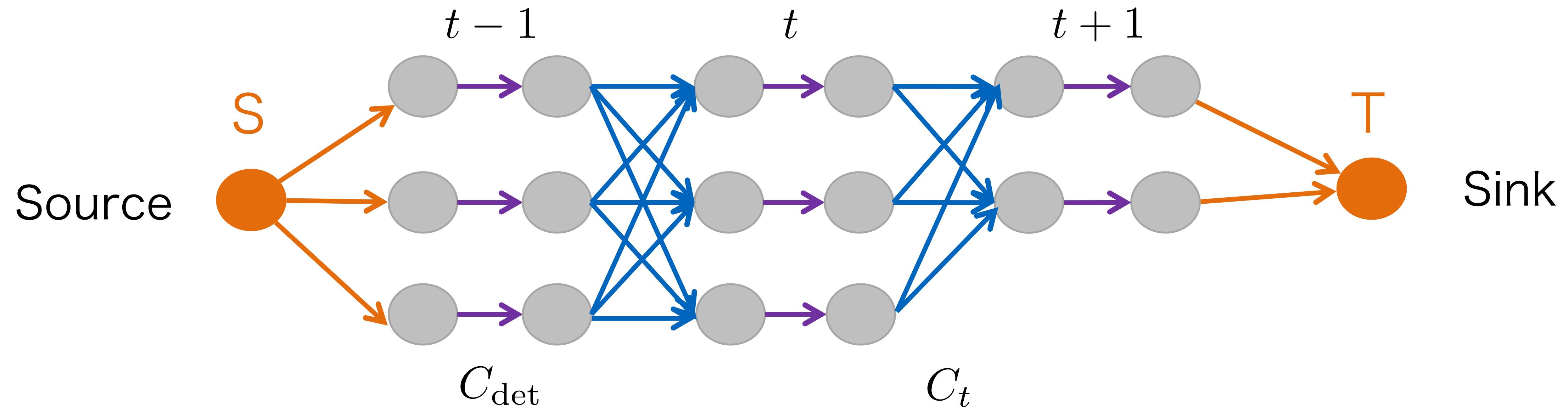
Zhang et al. “Global Data Association for Multi-Object Tracking Using Network Flows“. CVPR 2008

# Complete graph



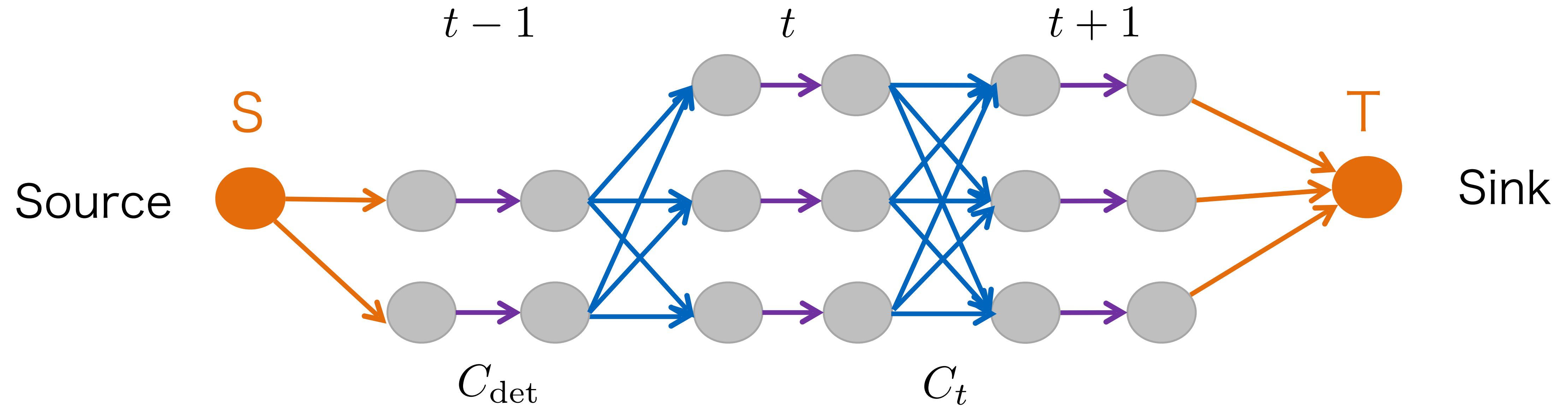
- The graph compactly encodes our problem.
  - QUIZ: How many trajectories (full-length) are possible?
  - But... detections are not perfect.

# Complete graph



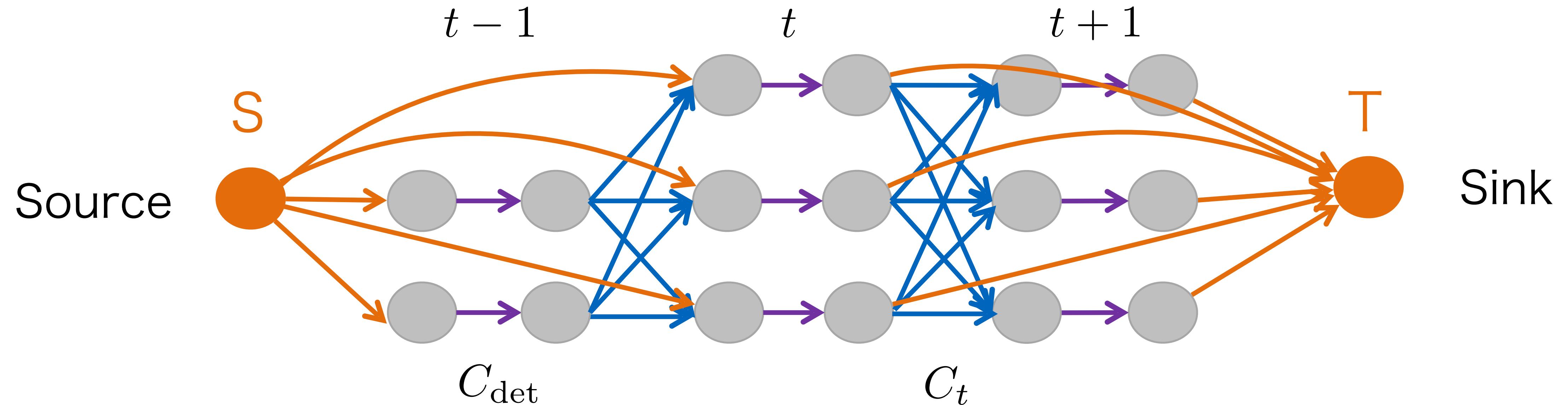
- If there is an occlusion, we find only two trajectories
  - e.g., occlusion in the last frame
- 'Y'

# Complete graph

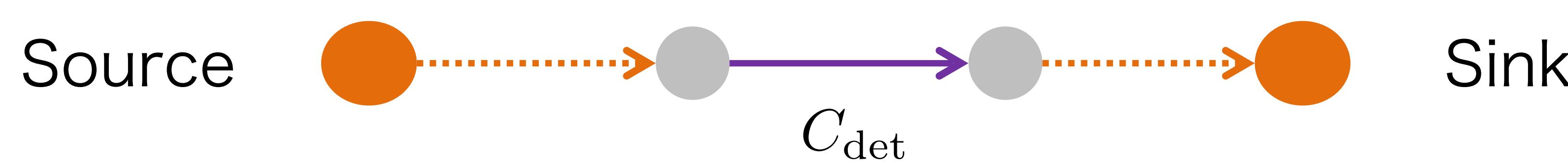


- If there is an occlusion, we find only two trajectories
  - or the object appears only in the second frame

# Complete graph



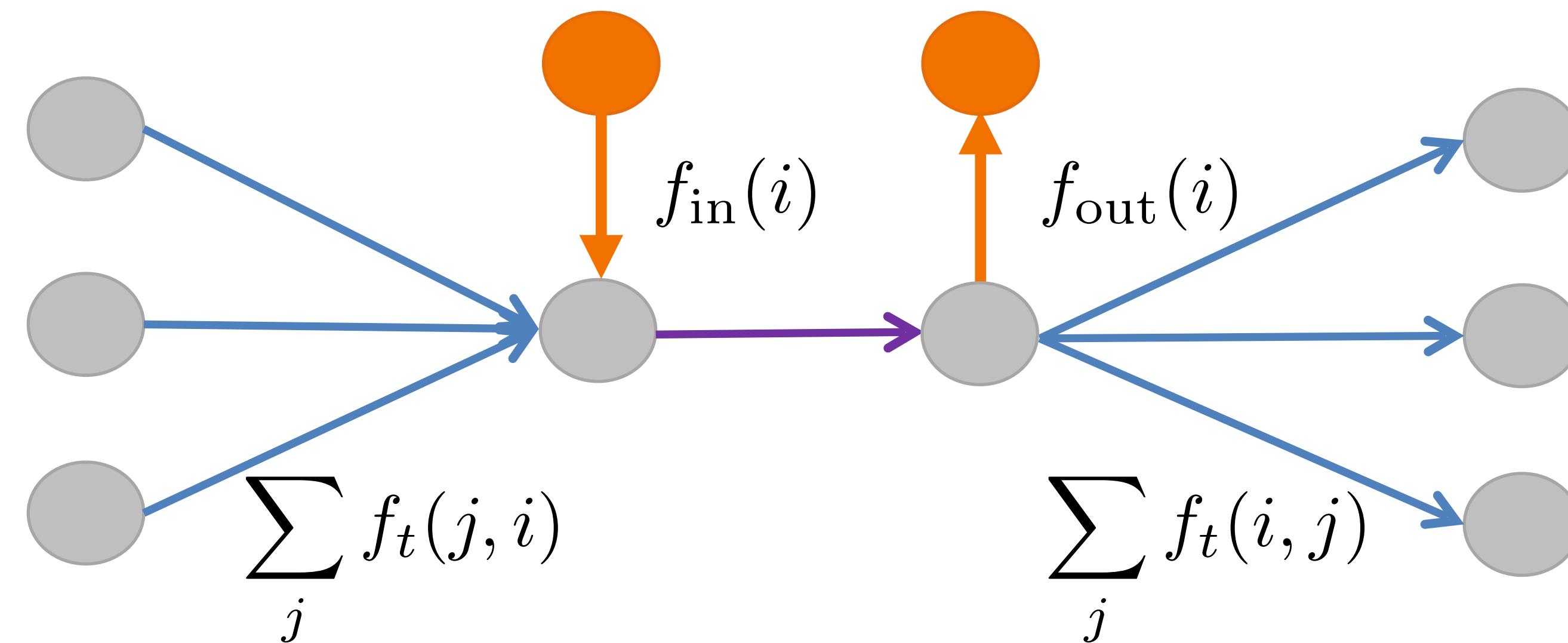
- Solution: Connect all nodes (detections) to entrance/exit nodes:



# Constraints

- Subject to **flow conservation**

$$\sum_j f_t(j, i) + f_{\text{in}}(i) = f(i) = \sum_j f_t(i, j) + f_{\text{out}}(i) \quad \forall f(\cdot, \cdot) \in \{0, 1\}$$
$$\forall f(\cdot) \in \{0, 1\}$$



# MAP formulation

$\rightarrow$  Maximum A posteriori

- Our solution is a set of trajectories  $\mathcal{T}^*$

$$\mathcal{T}^* = \arg \max_{\mathcal{T}} P(\mathcal{T} | \mathcal{X})$$

$$= \arg \max_{\mathcal{T}} P(\mathcal{X}|\mathcal{T})P(\mathcal{T})$$

Bayes rule

$$= \arg \max_{\mathcal{T}} \prod_i P(\mathbf{x}_i|\mathcal{T})P(\mathcal{T})$$

Assumption 1:  
Conditional independence of observations

$$= \arg \max_{\mathcal{T}} \prod_i P(\mathbf{x}_i|\mathcal{T}) \prod_{T_i \in \mathcal{T}} P(T_i)$$

Assumption 2:  
Independence of trajectories

See Leal-Taixé et al. (2011) for a more advanced model.

# MAP formulation

$$\arg \max_{\mathcal{T}} \prod_i P(\mathbf{x}_i | \mathcal{T}) \prod_{T_i \in \mathcal{T}} P(T_i)$$

$$= \arg \min_{\mathcal{T}} - \sum_i \log P(\mathbf{x}_i | \mathcal{T}) - \sum_{T_i \in \mathcal{T}} \log P(T_i)$$

log-space for optimisation

- ~~Trajectory model~~: count the entrance, exit and transition costs

$$P(T_i) = P_{\text{in}}(\mathbf{x}_0) \prod_{i=1, n} P_{\text{t}}(\mathbf{x}_i | \mathbf{x}_{i-1}) P_{\text{out}}(\mathbf{x}_n)$$

$$-\log P(T_i) = \log P_{\text{in}}(\mathbf{x}_0) - \sum_{i=1, n} \log P_{\text{t}}(\mathbf{x}_i | \mathbf{x}_{i-1}) - \log P_{\text{out}}(\mathbf{x}_n)$$

# Prior

$$\sum_{T_i \in \mathcal{T}} \log P(T_i)$$

- Trajectory model: count the entrance, exit and transition costs

$$T_i := \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$P(T_i) = P_{\text{in}}(\mathbf{x}_0) \prod_{j=1,n} P_t(\mathbf{x}_j \mid \mathbf{x}_{j-1}) P_{\text{out}}(\mathbf{x}_n)$$

how likely we reach  
this state from  
given state

$$-\log P(T_i) = -\log P_{\text{in}}(\mathbf{x}_0) - \sum_{j=1,n} \log P_t(\mathbf{x}_j \mid \mathbf{x}_{j-1}) - \log P_{\text{out}}(\mathbf{x}_n)$$

$$f_{\text{in}}(\mathbf{x}_0) C_{\text{in}}(\mathbf{x}_0) \quad f_t(\mathbf{x}_j, \mathbf{x}_i) C_t(\mathbf{x}_j, \mathbf{x}_i) \quad f_{\text{out}}(\mathbf{x}_n) C_{\text{out}}(\mathbf{x}_n)$$

# Likelihood

$$-\sum_i \log P(\mathbf{x}_i | \mathcal{T})$$

- We can use Bernoulli distribution:

$$P(\mathbf{x}_i | \mathcal{T}) := \begin{cases} \gamma_i, & \text{if } \exists T_j \in \mathcal{T}, \mathbf{x}_i \in T_j \\ 1 - \gamma_i, & \text{otherwise} \end{cases}$$

- $\gamma_i$  denotes prediction confidence (e.g. provided by the detector)

$$\begin{aligned} -\log P(\mathbf{x}_i | \mathcal{T}) &= -f(\mathbf{x}_i) \log \gamma_i - (1 - f(\mathbf{x}_i)) \log(1 - \gamma_i) \\ &= f(\mathbf{x}_i) \log \frac{1 - \gamma_i}{\gamma_i} - \log(1 - \gamma_i) \\ &\quad \left| \begin{array}{l} C_{\text{det}}(\mathbf{x}_i) \\ \hline \end{array} \right. \end{aligned}$$

can be ignored  
in optimisation

# Optimisation

- $(C_{\text{det}}, C_{\text{in}}, C_{\text{out}}, C_t)$  are estimated from data. Then:

- Construct the graph  $G(V, E, C, f)$  from observation set  $\mathcal{X}$
- Start with empty flow
- WHILE (  $f(G)$  can be augmented )

- Augment  $f(G)$  by one.
  - Find the min cost flow by the algorithm of [12].
  - IF ( current min cost < global optimal cost )  
Store current min-cost assignment as global optimum.

- Return the global optimal flow as the best association hypothesis

Binary/Fibonacci search  
 $O(\log n)$

Min-cost flow algorithm  
 $O(n^2 m \log n)$

# Summary



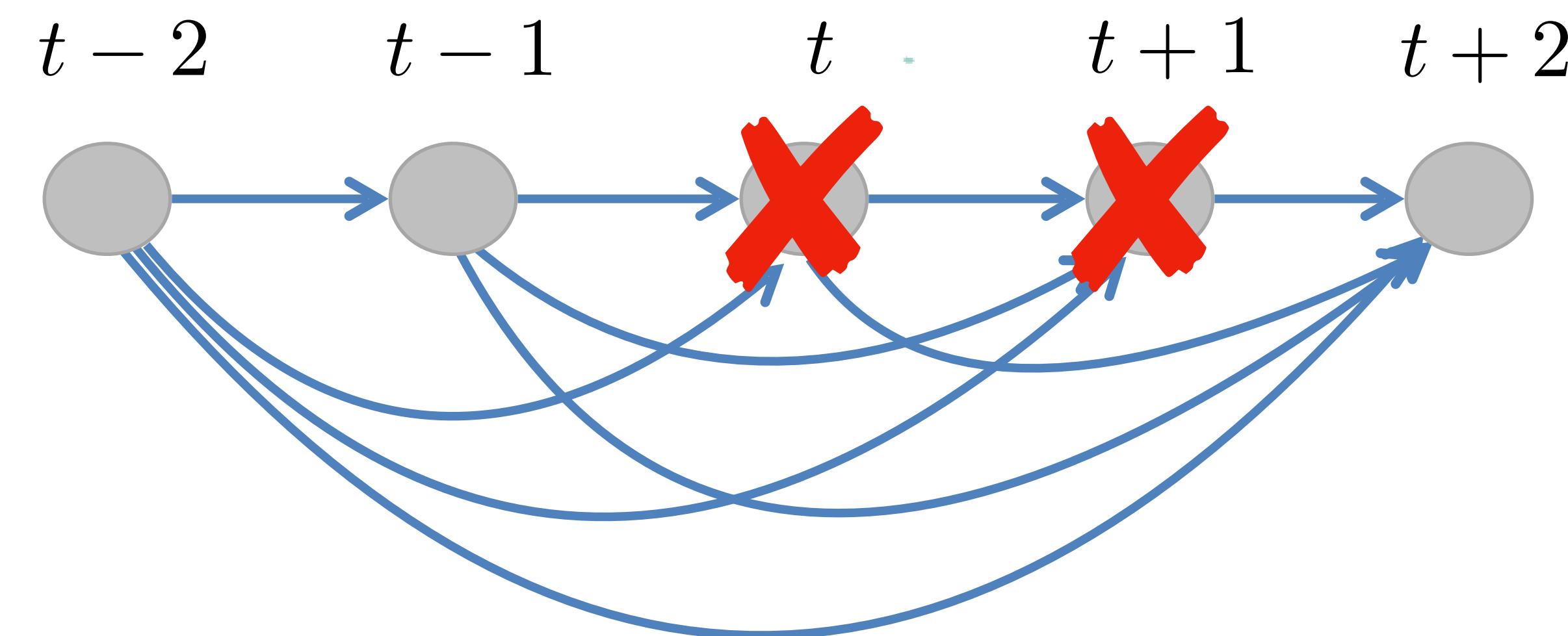
- Min-cost max-flow formulation:
  - the maximum number of trajectories with minimum costs.
- Optimisation maximises maximum a-posteriori (MAP):
  - global solution and efficient (polynomial time).
- Open questions:
  - how to handle occlusions?

# Open questions

- How to handle occlusions?

# Handling occlusions

- How to handle occlusions?
    - “Markovian” formulation may not be sufficient ( $\rightarrow$  dense graphs)
- $\rightarrow$  what if 1 step is not exist.



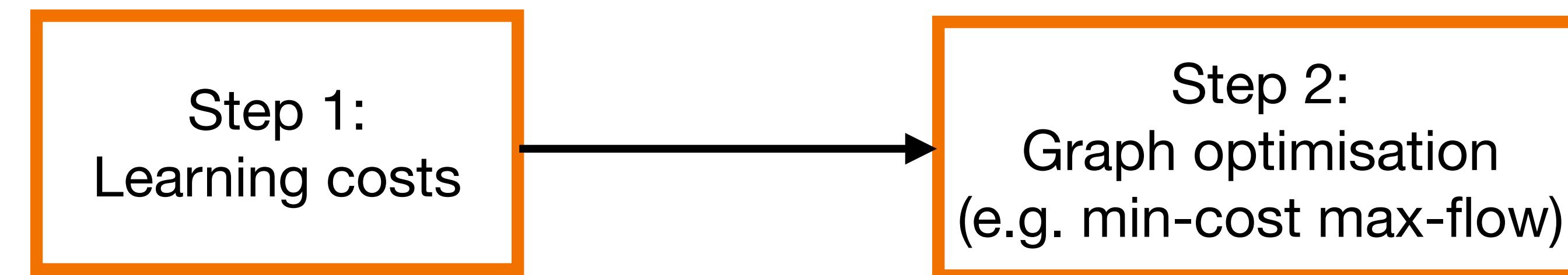
# Handling occlusions

More advanced (dense) graph formulations:

- M. Keuper et al. „Motion segmentation and multiple object tracking by correlation co-clustering“. PAMI 2018.
- S. Tang et al. „Subgraph decomposition for multi-target tracking: CVPR 2015.
- S. Tang et al. „Multiple people tracking by lifted multicut and person reidentification“. CVPR 2017

# Open questions

- How to handle occlusions?
- **How to learn costs?**

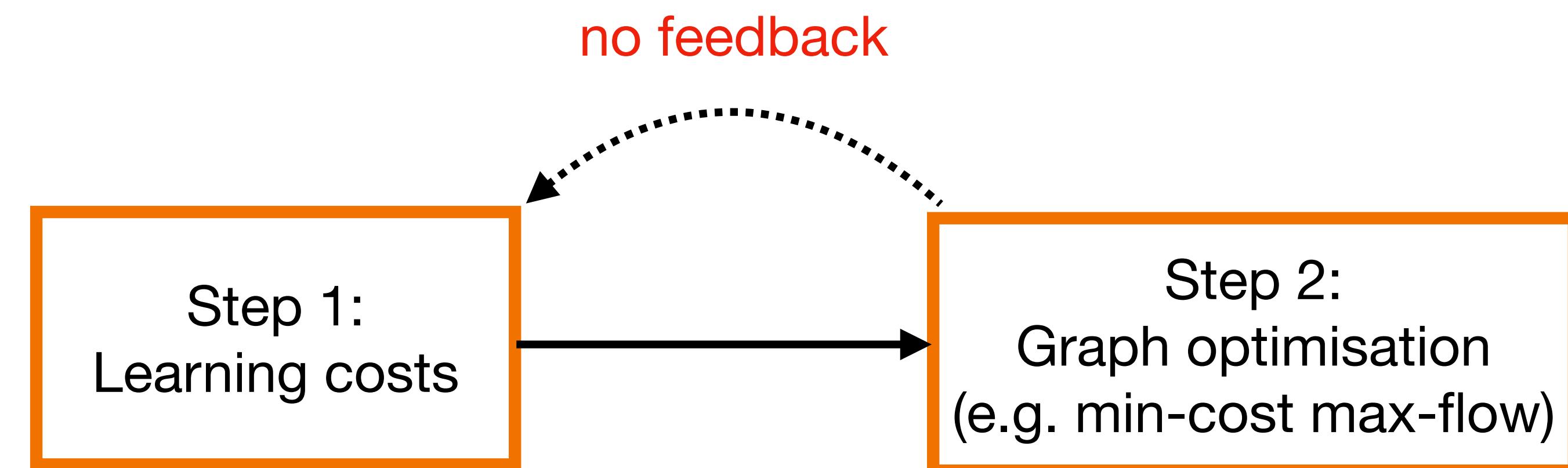


# Improving costs

- L. Leal-Taixé et al. “Learning an image-based motion context for multiple people tracking”. CVPR 2014.
- L. Leal-Taixé et al. “Everybody needs somebody: Modeling social and grouping behavior on a linear programming multiple people tracker“. ICCVW 2011
- L. Leal-Taixé et al. “Learning by tracking: Siamese CNN for robust target association”. CVPRW 2016.
- S. Schulter et al. „Deep network flow for multi-object tracking“. CVPR 2017.
- J. Son at al. „Multi-object tracking with quadruplet convolutional neural networks“. CVPR 2017.
- Ristani and Tomasi. “Features for multi-target multi-camera tracking and re-identification”. CVPR 2018
- J. Xu et al. „Spatial-temporal relation networks for multi-object tracking“. ICCV 2019.

# Open questions

- How to handle occlusions?
- **How to learn costs?**
  - costs may be specific to the graph formulation and optimisation



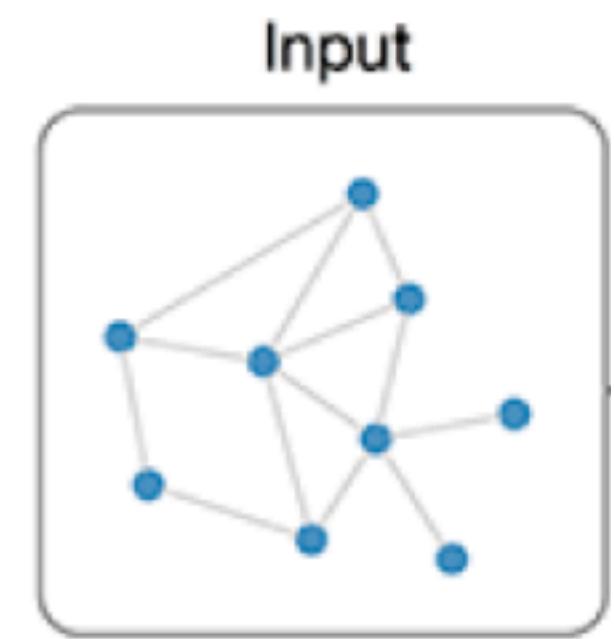
# End-to-end learning?

- Can we learn:
  - features for multi-object tracking (encoding costs);
  - to find solution on the graph?
- Goal: Generalise the graph structure we have used and perform end-to-end learning

# Message Passing Networks

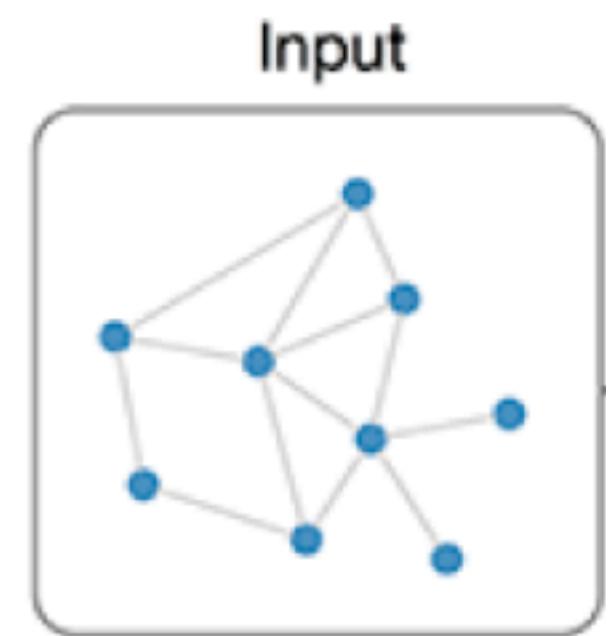
# Setup

- Input: task-encoding graph
  - nodes: detections encoded as feature vectors
  - edges: node interaction (e.g. inter-frame)
- Output: graph partitioning into (disjoint) trajectories
  - e.g. encoded by edge label (0,1)

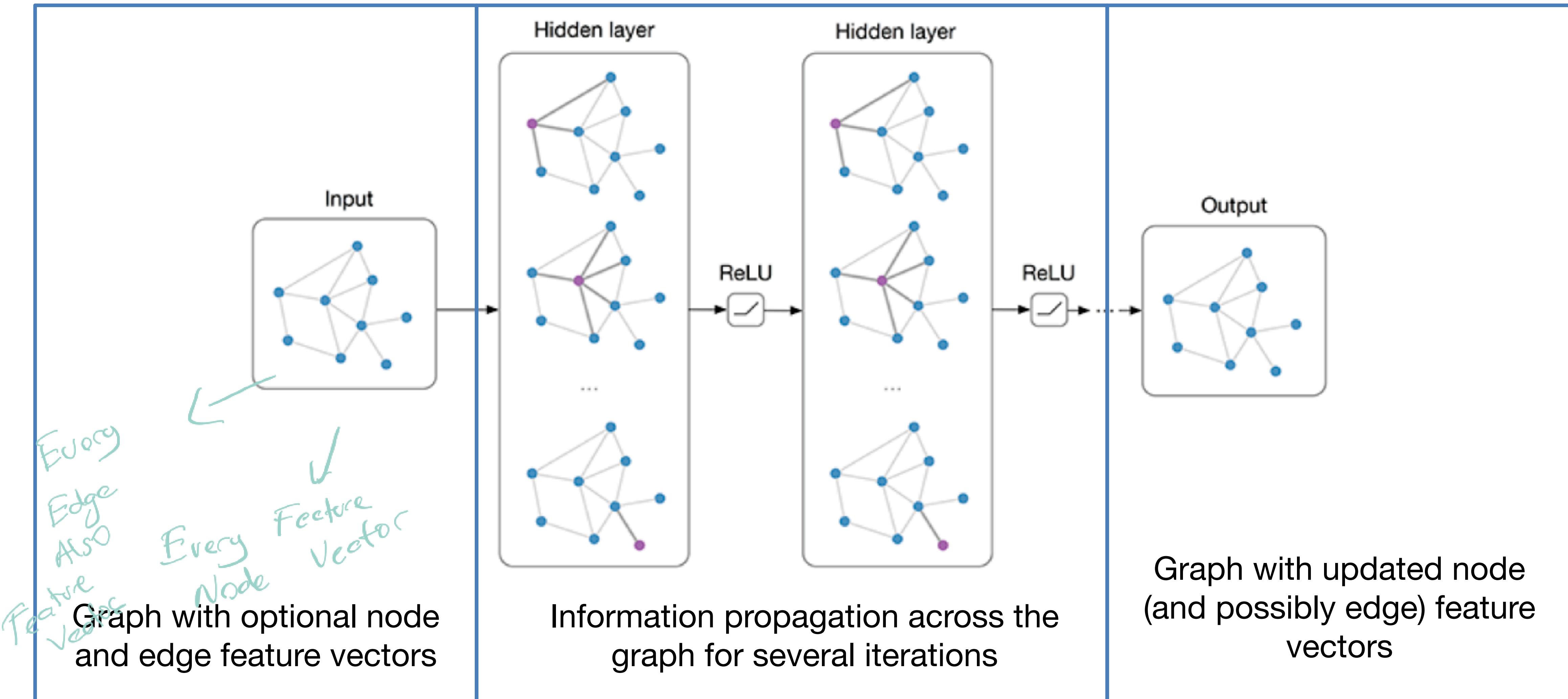


# Deep learning on graphs

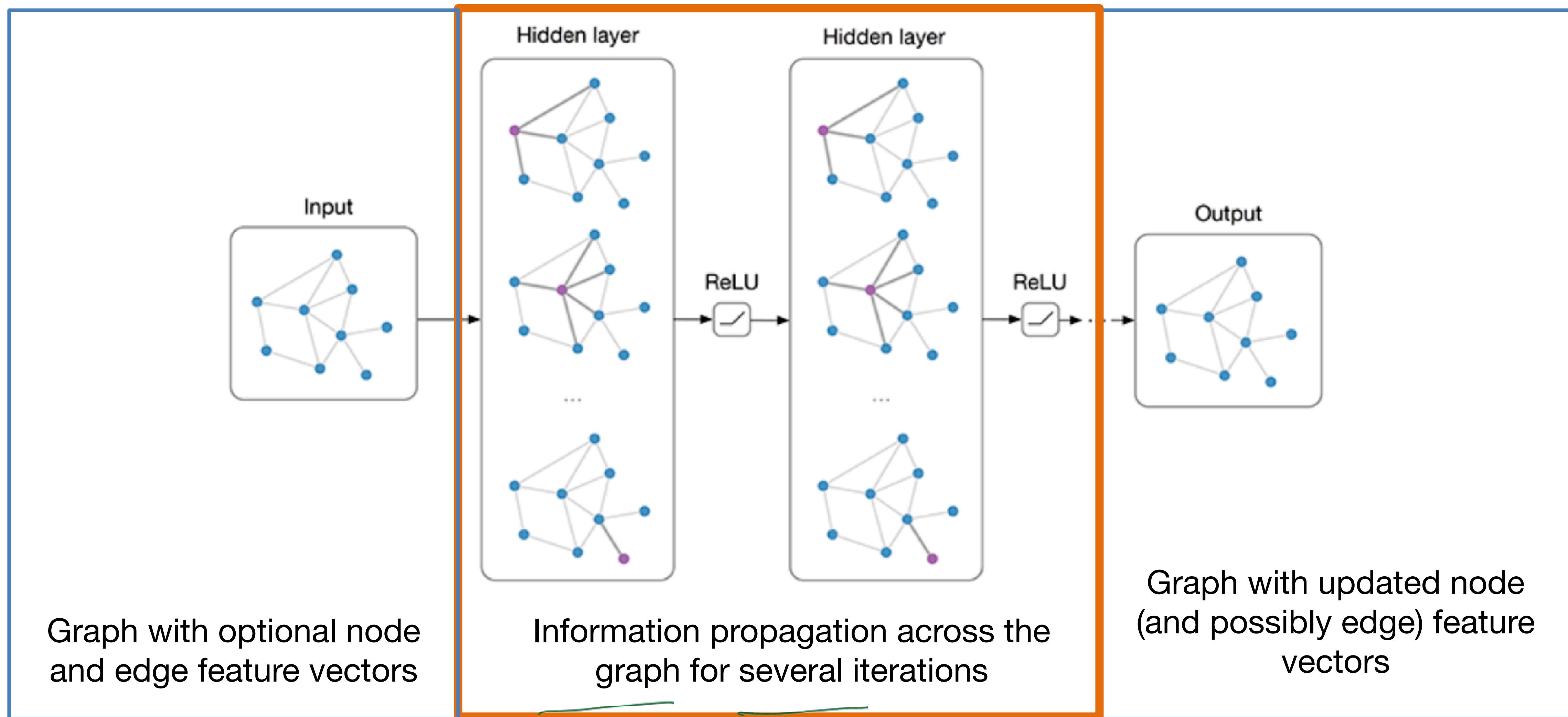
- Key challenges:
  - Graph can be of arbitrary size  
(number of nodes and edges)
  - Need invariance to node permutations



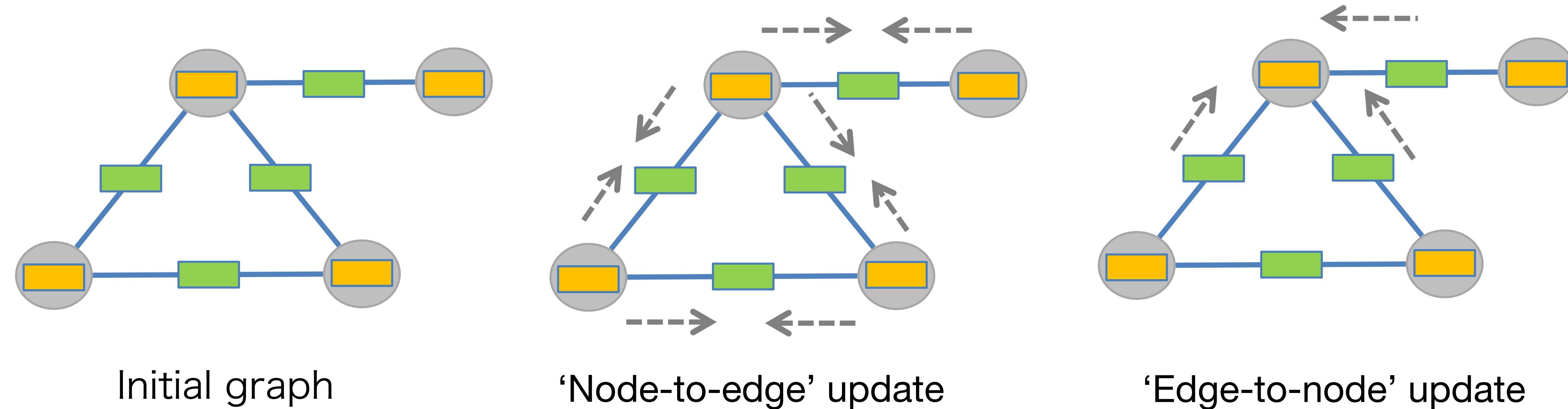
# Message Passing Networks



# Message Passing Networks



# Message passing



- We can divide the propagation process in two steps:

- 'node-to-edge' and 'edge-to-node' updates.

- Alternate these two updates (message passing)

→ encodes context into embeddings.

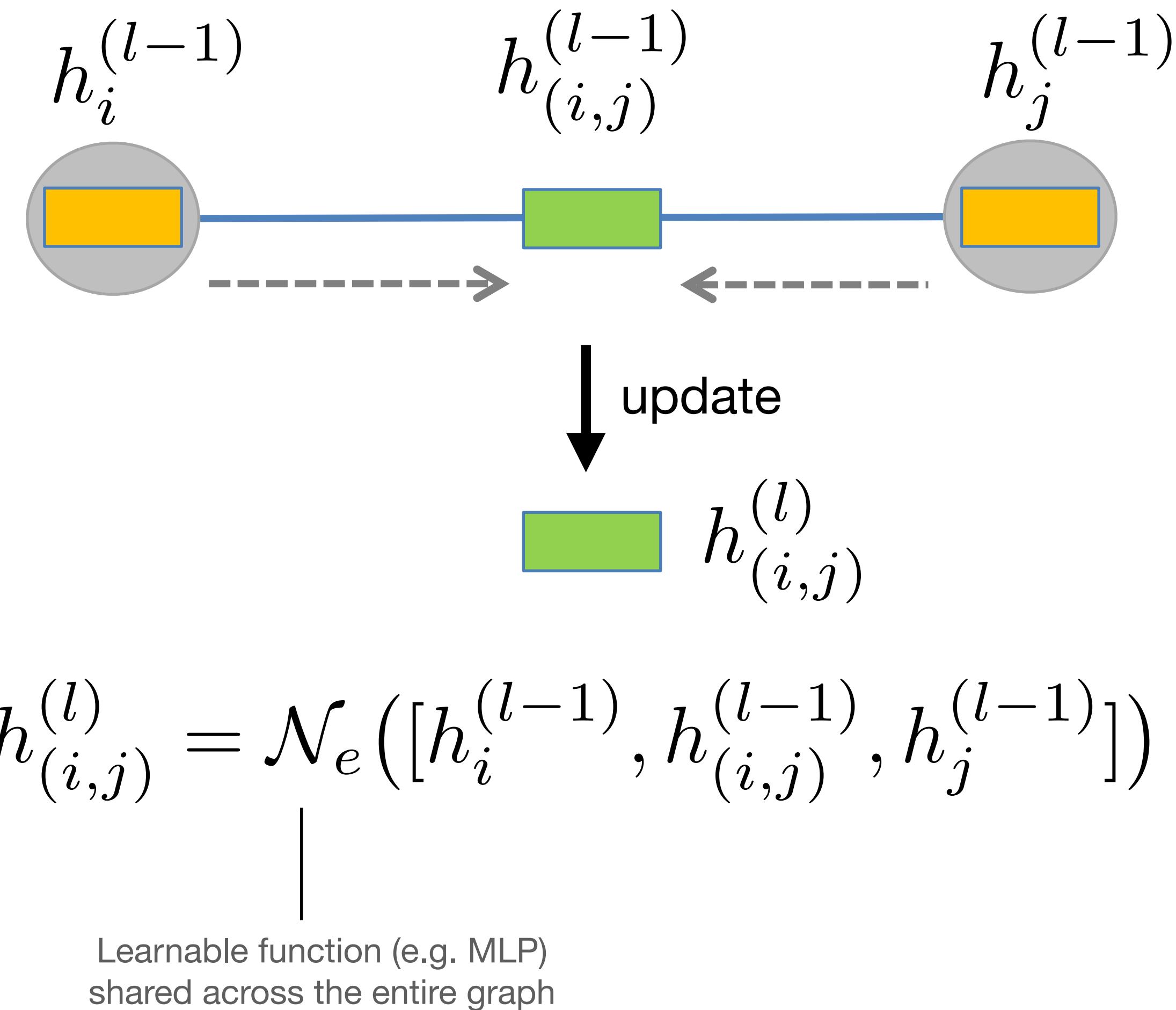
Node embeddings

Edge embeddings

# Notation

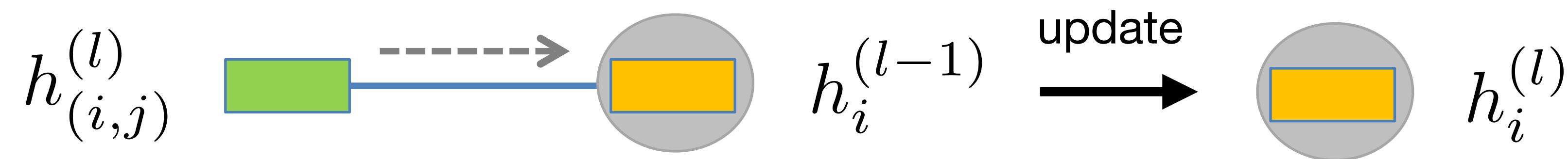
- Graph:  $G = (V, E)$
- Initial embeddings:
  - Node embedding:  $h_i^{(0)}, i \in V$
  - Edge embedding:  $h_{(i,j)}^{(0)}, (i, j) \in E$
- Embeddings after  $l$  steps:  $h_i^{(l)}, i \in V$        $h_{(i,j)}^{(l)}, (i, j) \in E$

# Node-to-edge update



# Edge-to-node updates

- Use the updated edge embeddings to update nodes:



- After a round of edge updates, each edge embedding contains information about its pair of incident nodes.

- By analogy:  $h_i^{(l)} = \mathcal{N}_v([h_i^{(l-1)}, h_{(i,j)}^l])$
- In general, we may have an arbitrary number neighbours (“degree”, or “valency”)

# Edge-to-node updates

- Define a permutation-invariant aggregation function:

$$\Phi^{(l)}(i) := \Phi\left(\left\{h^{(l)}(i, j)\right\}_{j \in Ne(i)}\right)$$

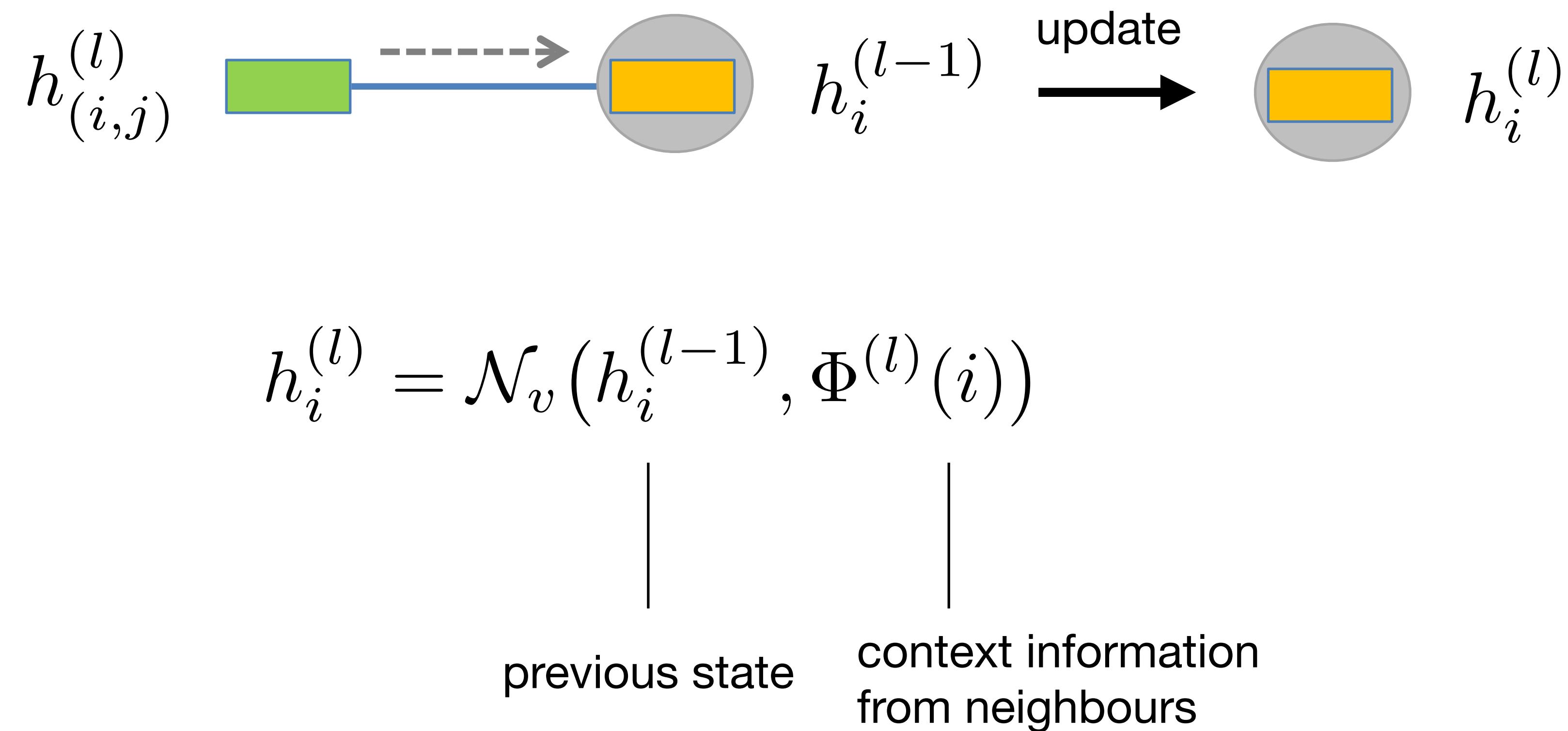
The input is a **set** of embeddings  
from incident edges

- What are examples ~~permutation-invariant~~ functions? (QUIZ)
- sum/mean, max/min (commutative operations).

Qi et al., "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation" (2017)

# Edge-to-node updates

- Re-define the edge-to-node updates for a general graph:



# Remarks

- Main goal: gather context information into node and edge embeddings
- Is one iteration of node-to-edge/edge-to-node updates enough? (QUIZ)
- One iteration increases the receptive field of a node/edge by 1  $\rightarrow N_0$ 
  - in practice: iterate message passing multiple times (hyperparameter).
- All operations used are differentiable!
- All vertices/edges are treated equally, i.e. the parameters are shared.

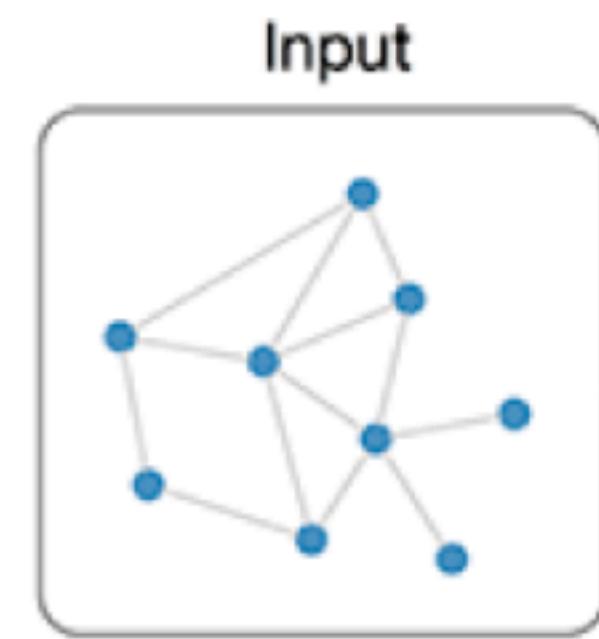
# Deep learning on graphs

- Several works have proposed generalizations of Neural Networks that can operate on graph-structured domains:
  - Scarselli et al. “The Graph Neural Network Model”. IEEE Trans. Neur. Net 2009.
  - Kipf et al. “Semi-Supervised Classification with Graph Convolutional Networks”. ICLR 2016.
  - Gilmer et al. “Neural Message Passing for Quantum Chemistry”. ICML 2017
  - Battaglia et al. “Relational inductive biases, deep learning, and graph networks”. arXiv 2018 (review paper)
- Key challenges:
  - Variable sized inputs (number of nodes and edges)
  - Need invariance to node permutations

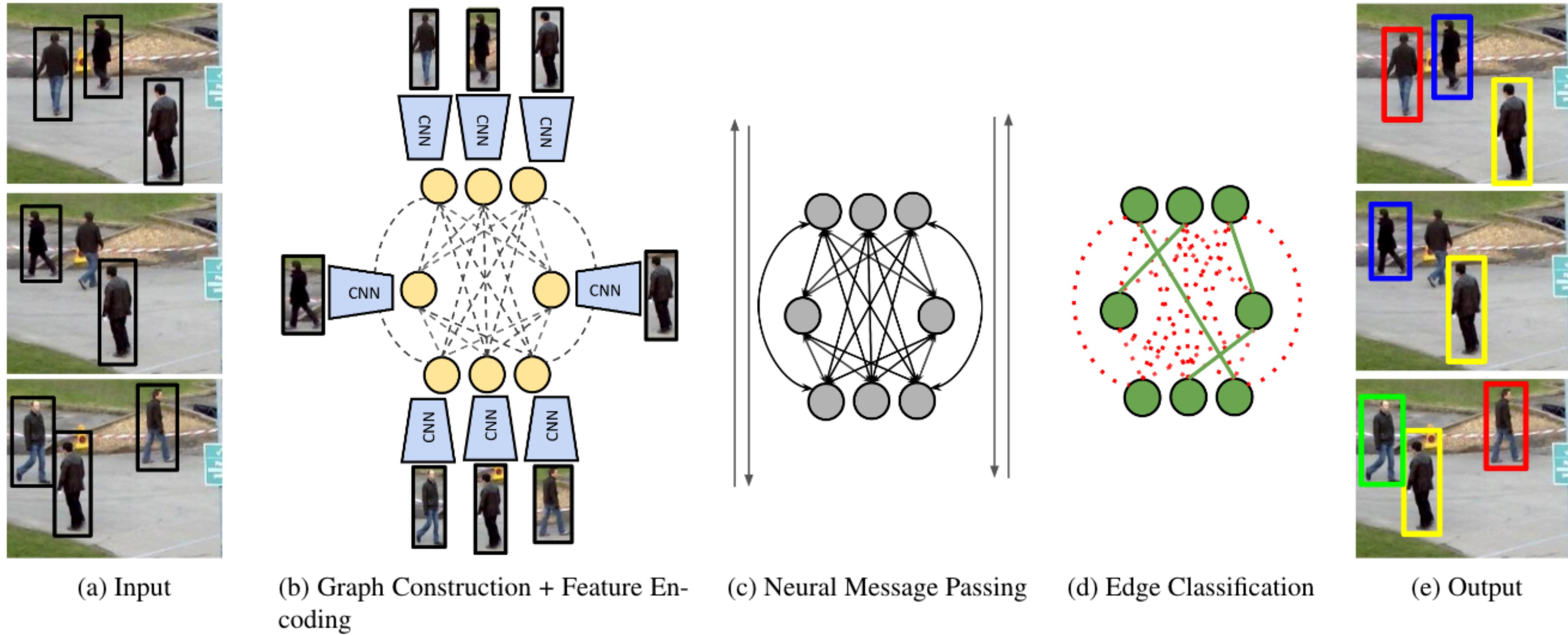
# MOT with Message Passing Networks

# Recall our setup

- Input: task-encoding graph
  - nodes: detections encoded as feature vectors
  - edges: node interaction (e.g. inter-frame)
- Output: graph partitioning into (disjoint) trajectories
  - e.g. encoded by edge label (0,1)



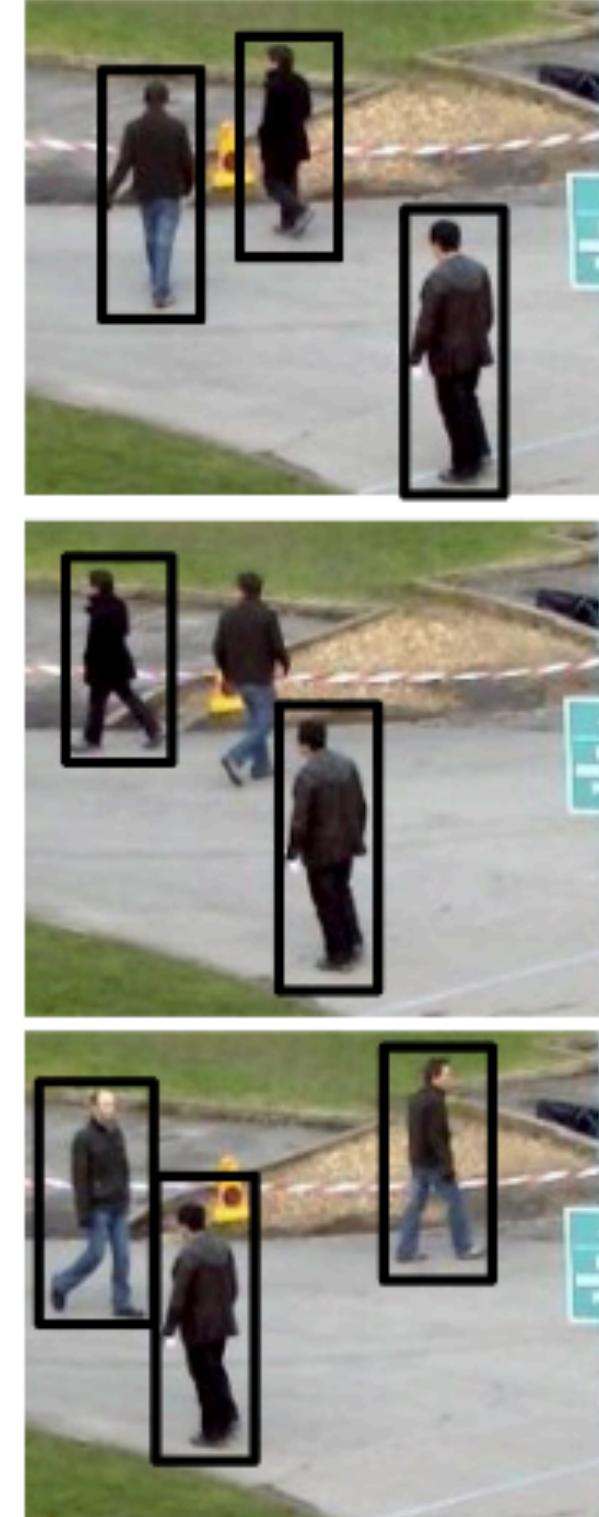
# Overview



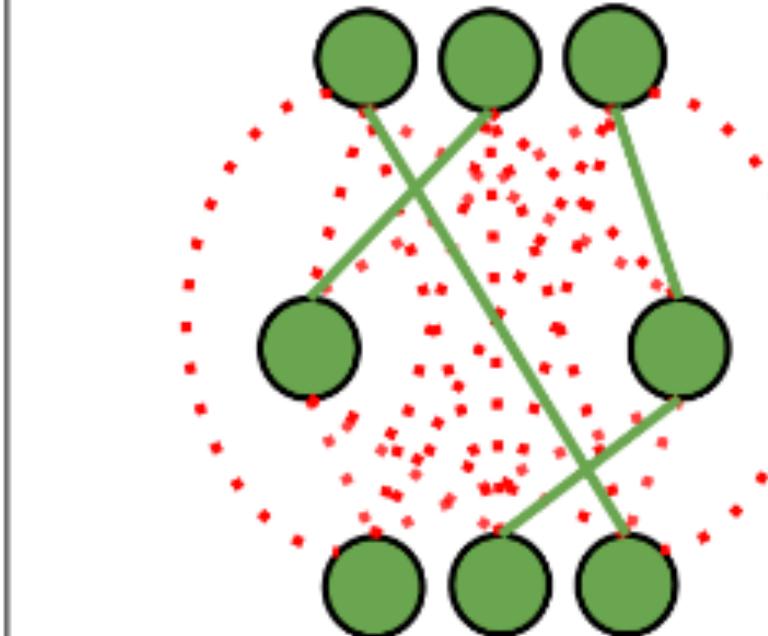
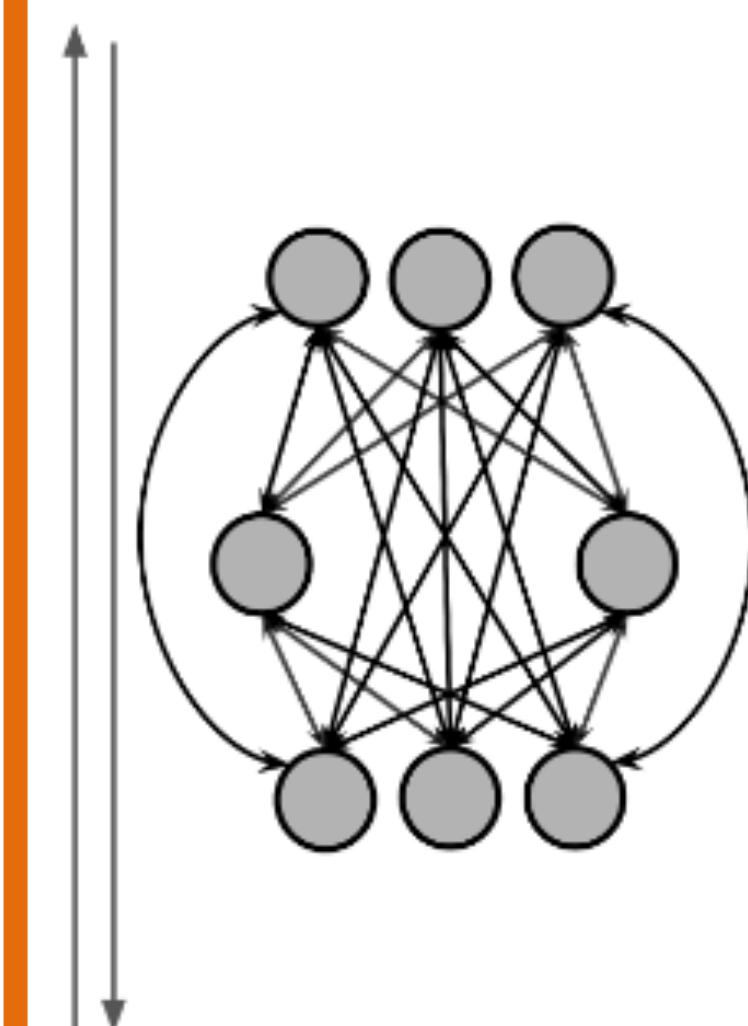
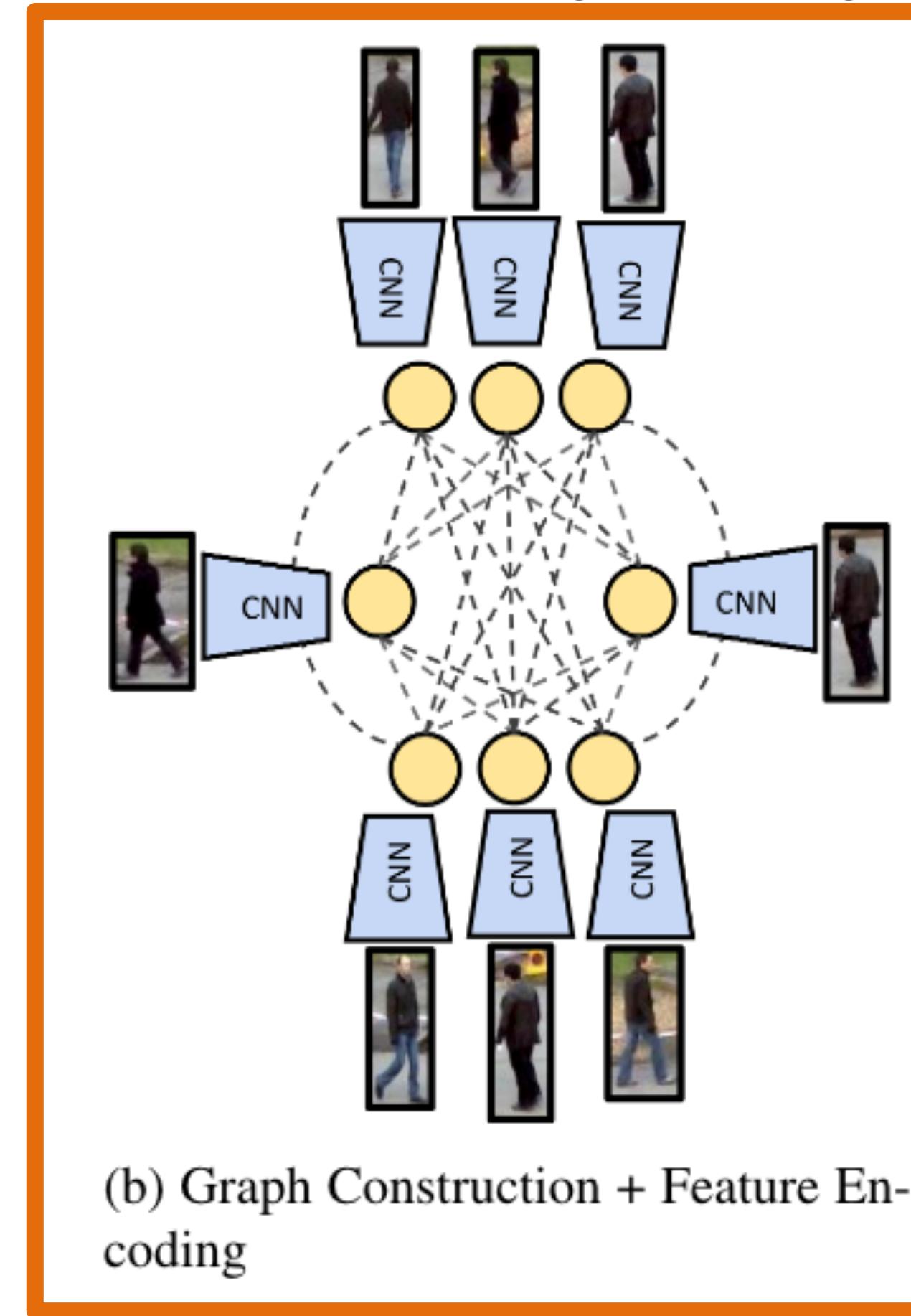
Brasó and Leal-Taixé. "Learning a Neural Solver for Multiple Object Tracking" (2019).

# Overview

Encode appearance and scene geometry cues into node and edge embeddings



(a) Input



(e) Output

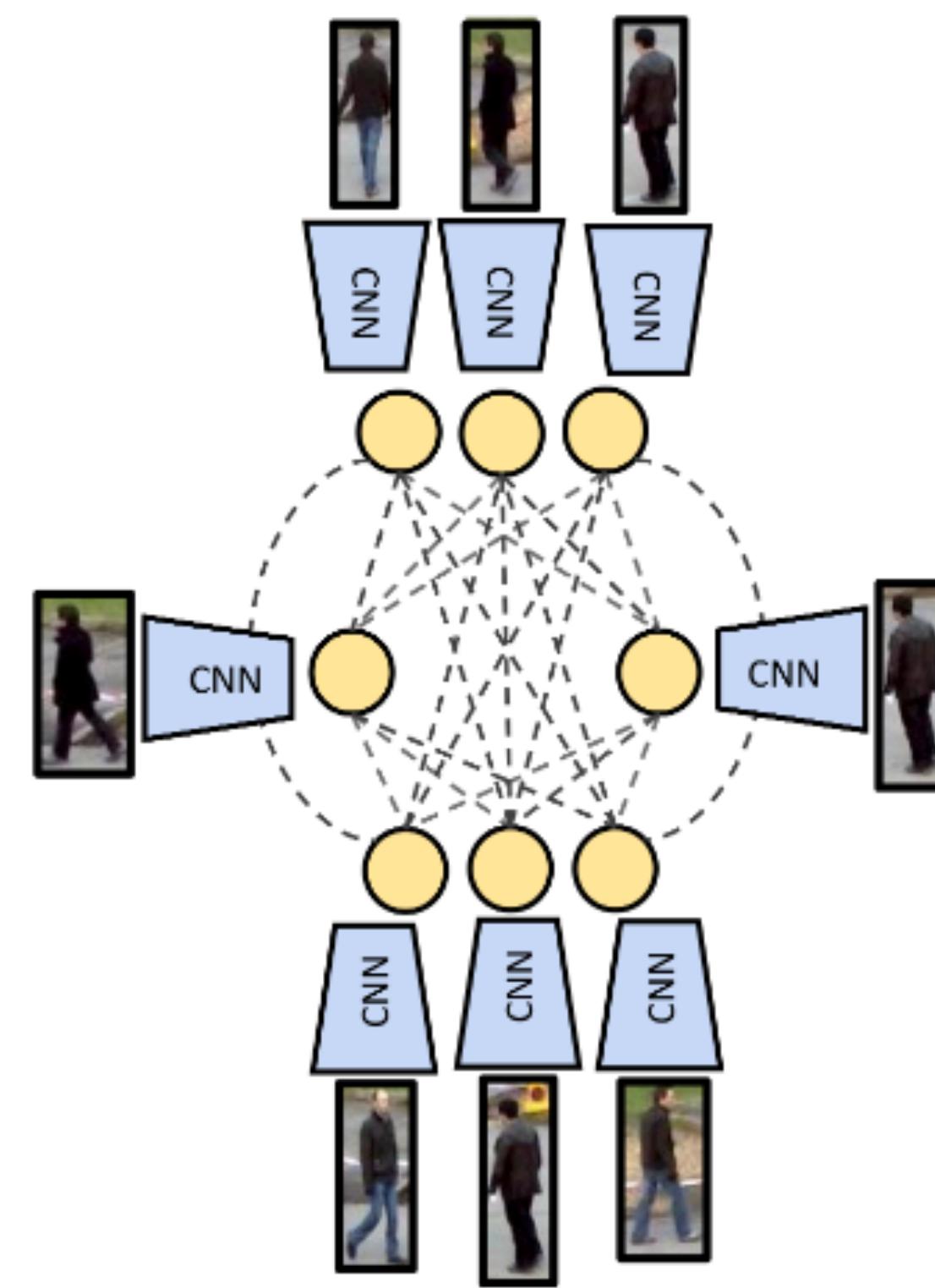
Brasó and Leal-Taixé. "Learning a Neural Solver for Multiple Object Tracking" (2019).

# Overview

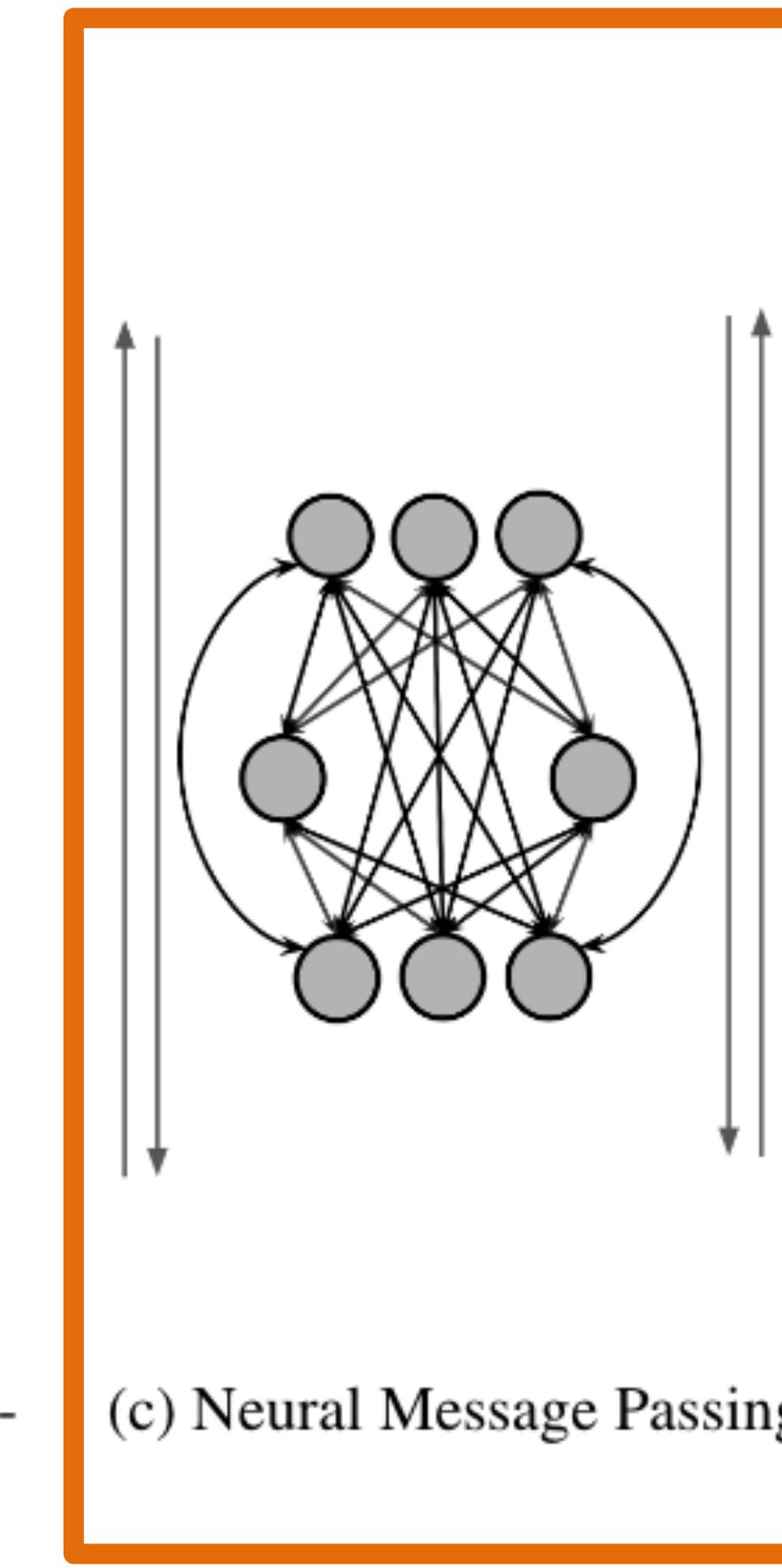
Propagate cues across the entire graph  
with neural message passing



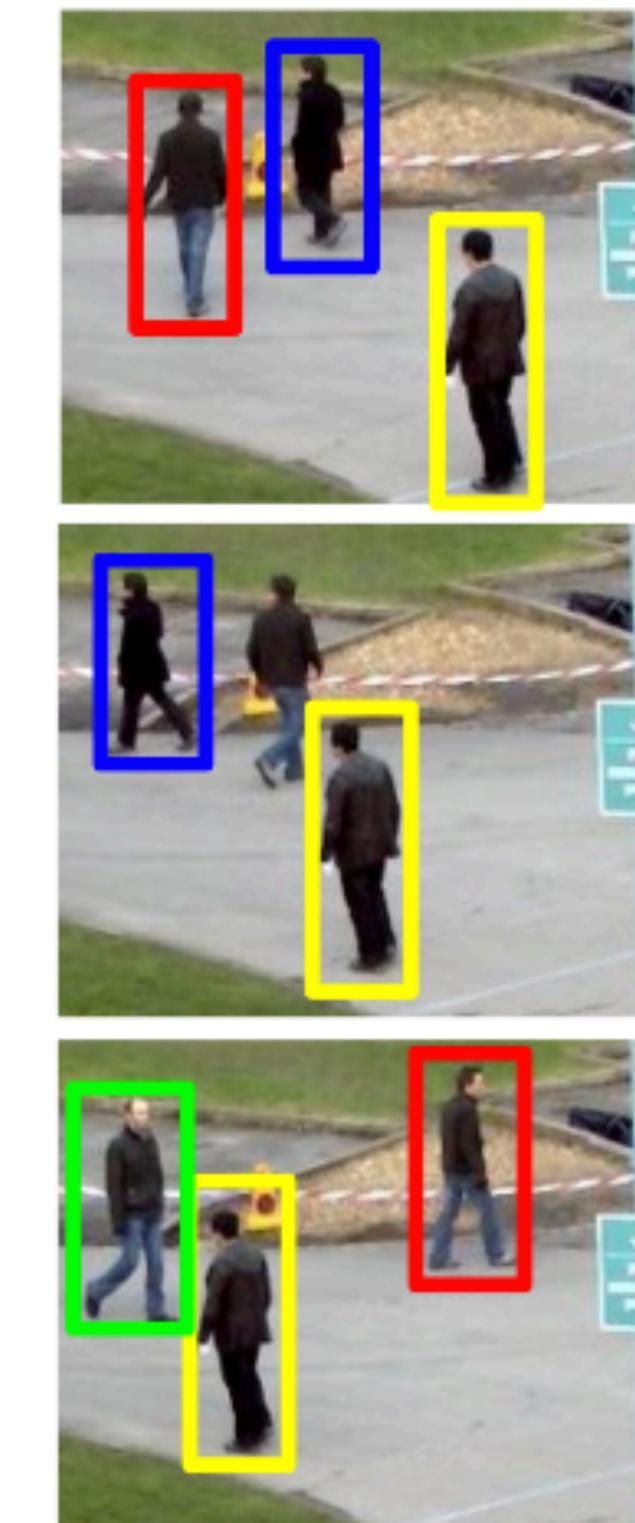
(a) Input



(b) Graph Construction + Feature Encoding



(c) Neural Message Passing



(d) Edge Classification

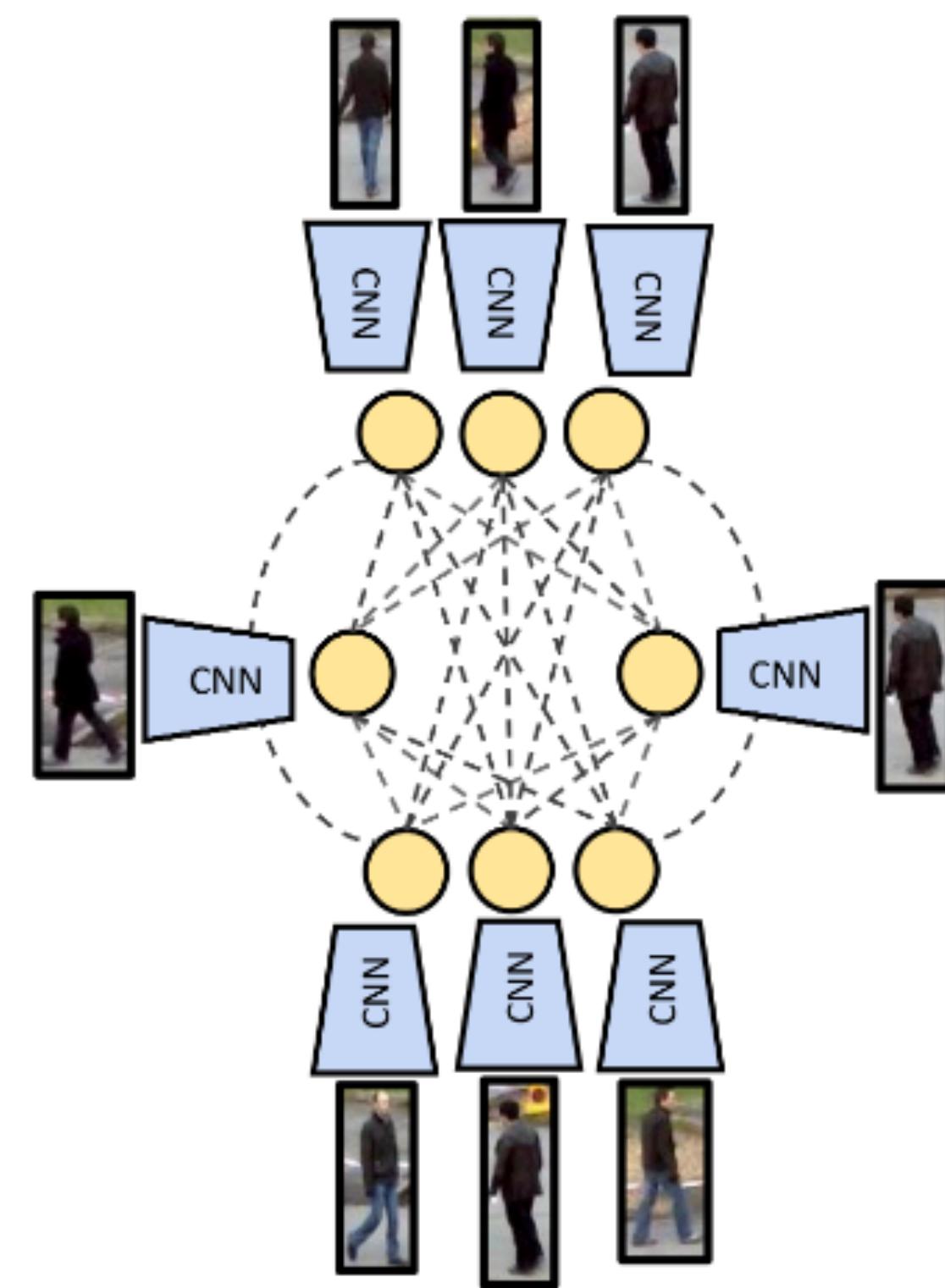
Brasó and Leal-Taixé. "Learning a Neural Solver for Multiple Object Tracking" (2019).

# Overview

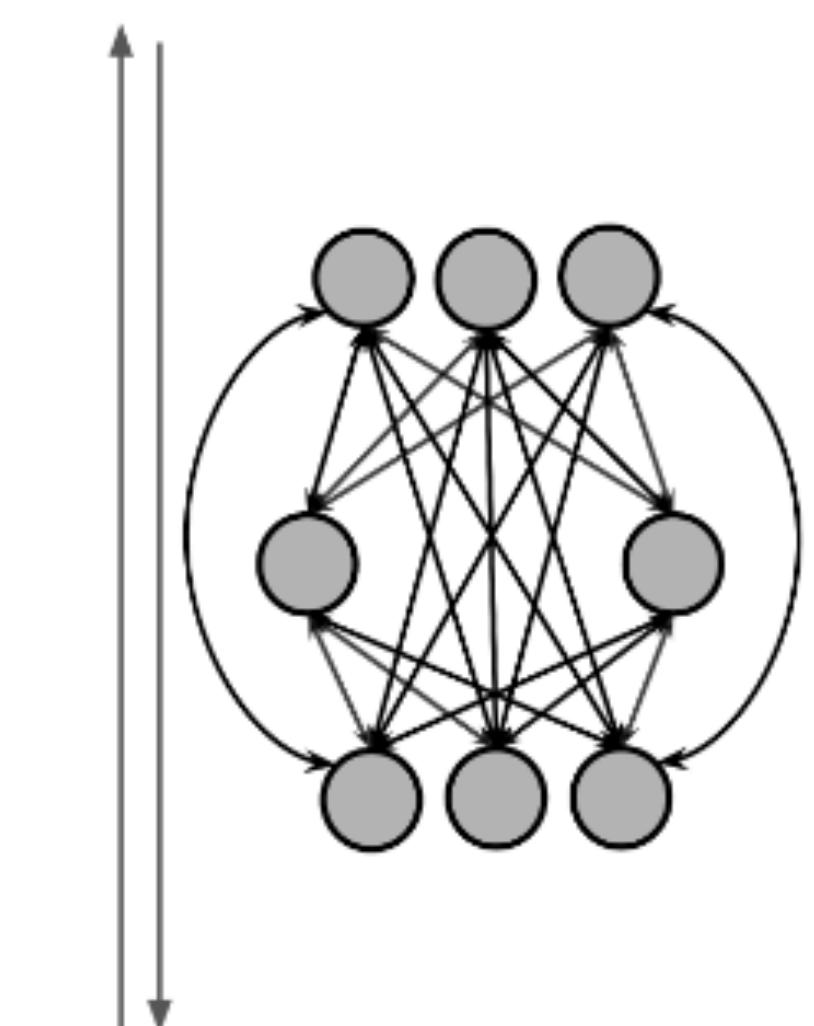
Learn to directly predict solutions  
of the Min-Cost Flow problem by  
classifying edge embeddings



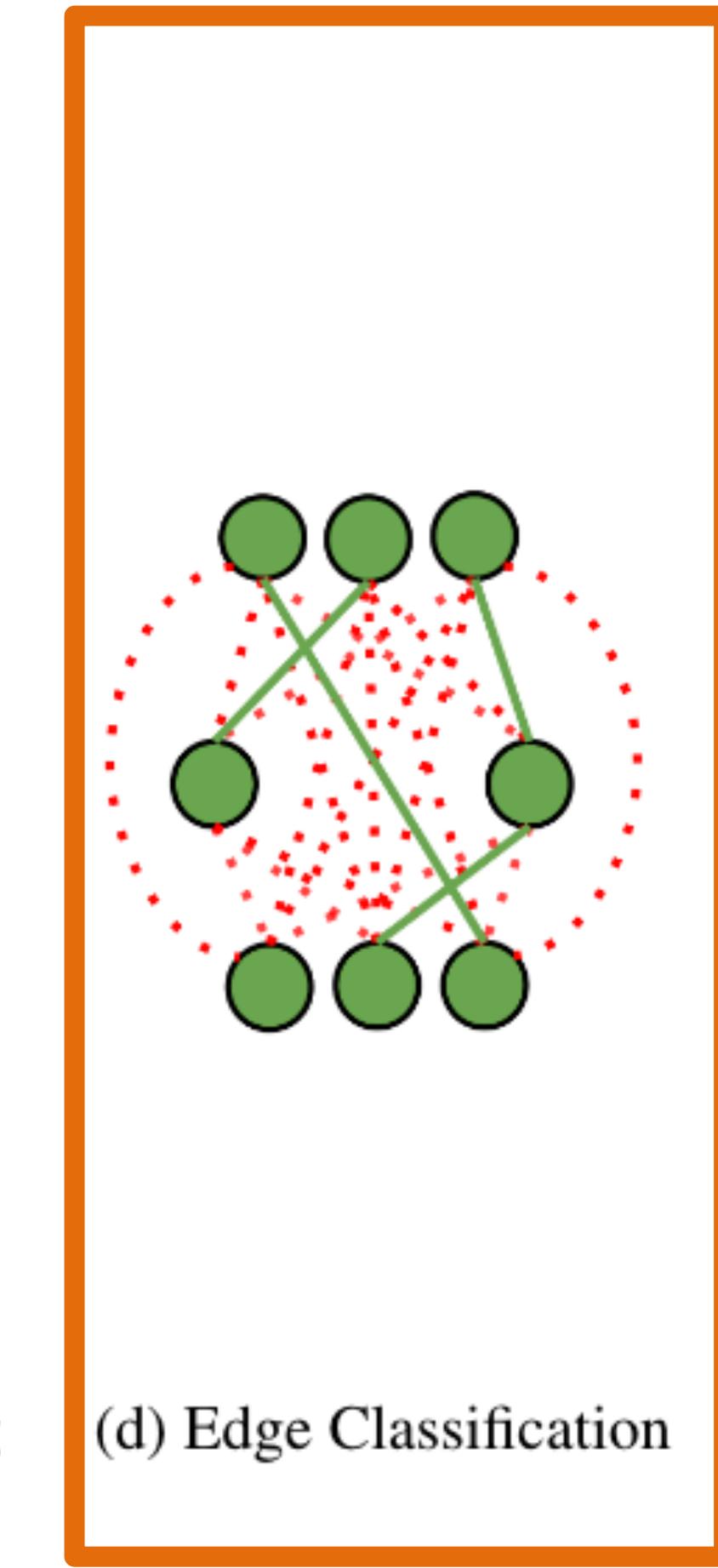
(a) Input



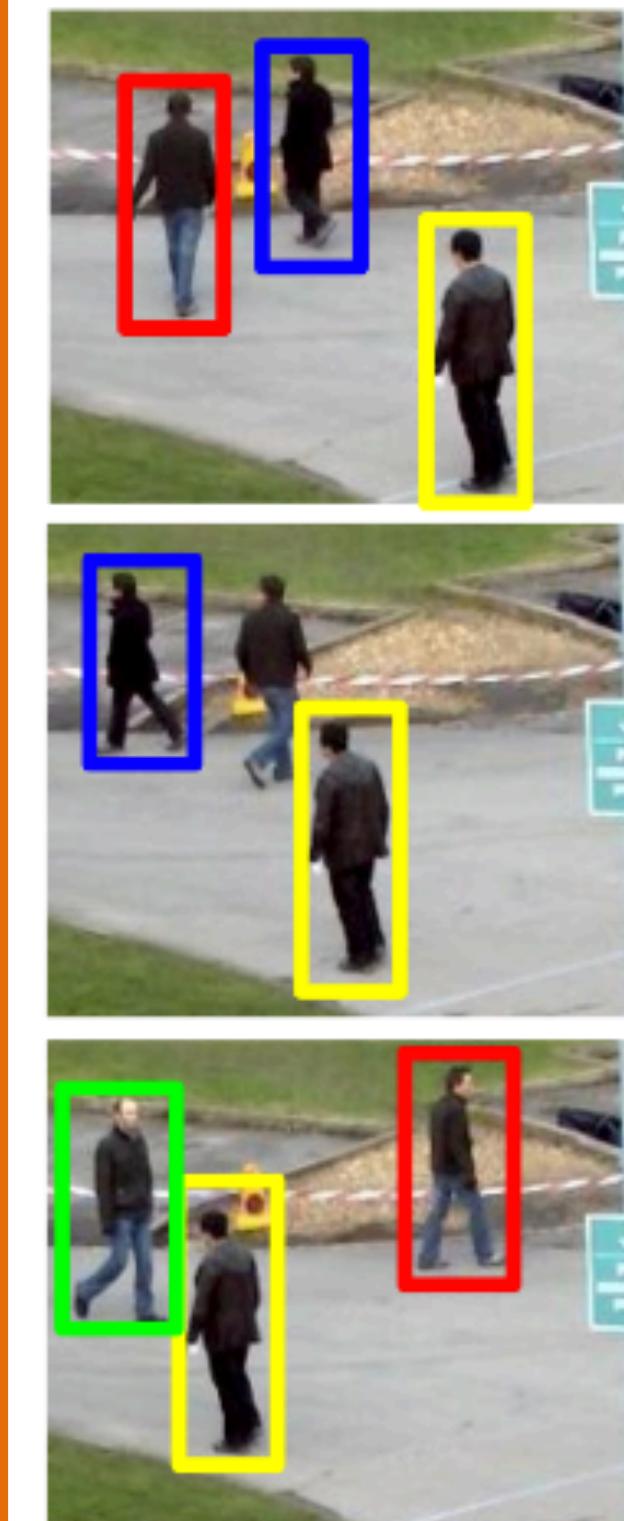
(b) Graph Construction + Feature Encoding



(c) Neural Message Passing



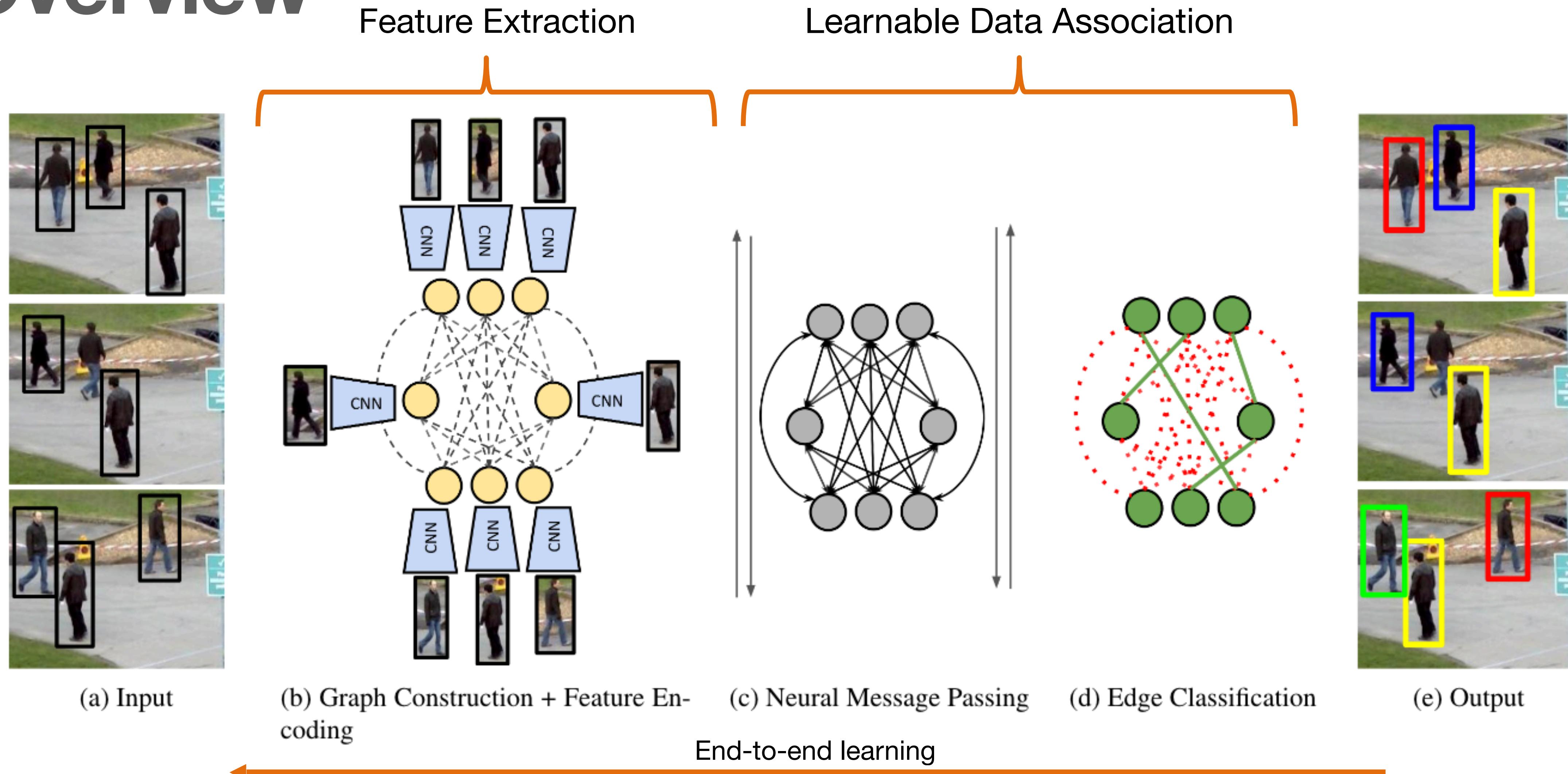
(d) Edge Classification



(e) Output

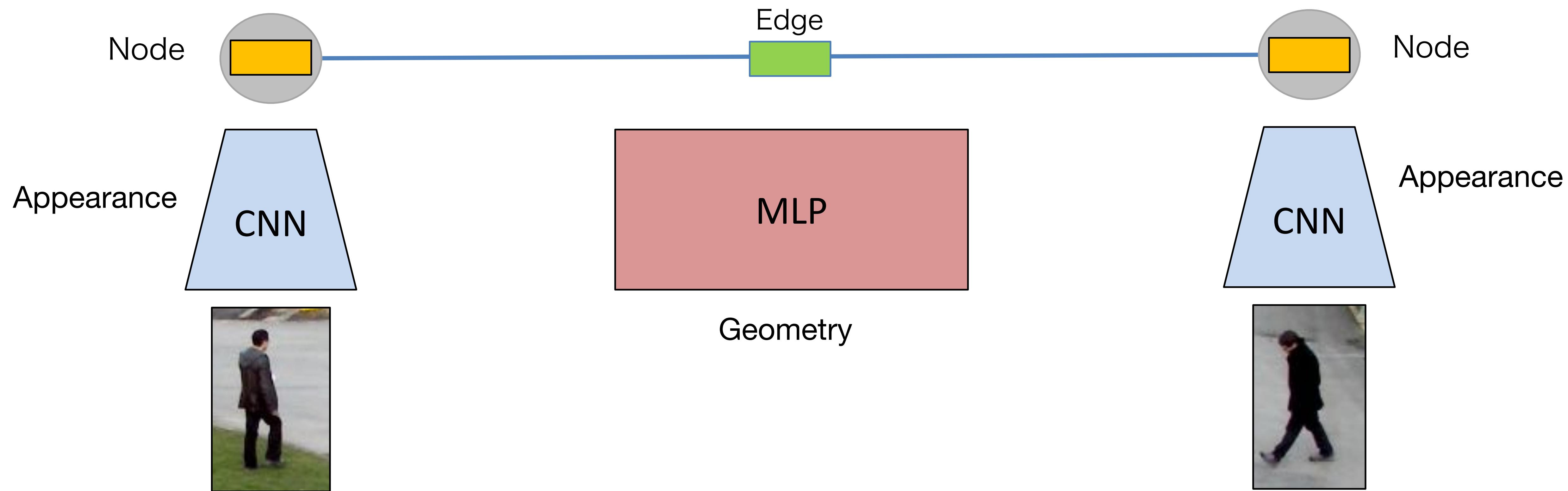
Brasó and Leal-Taixé. "Learning a Neural Solver for Multiple Object Tracking" (2019).

# Overview



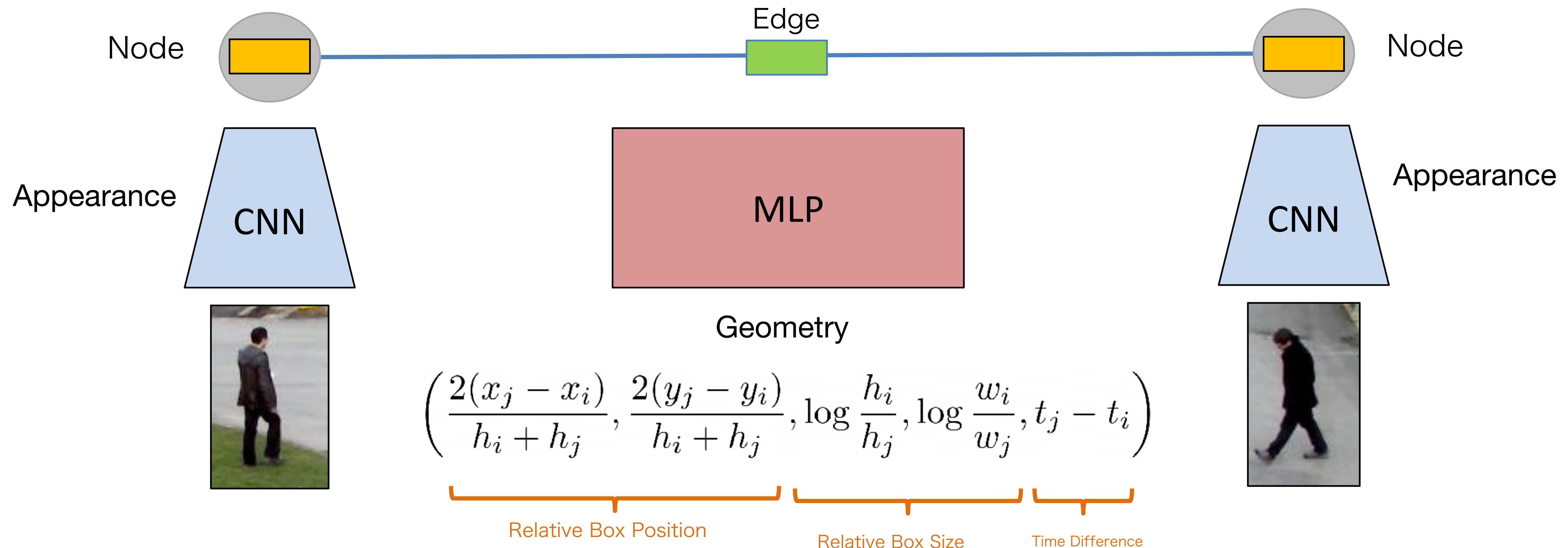
# Feature encoding

- Appearance and geometry encodings



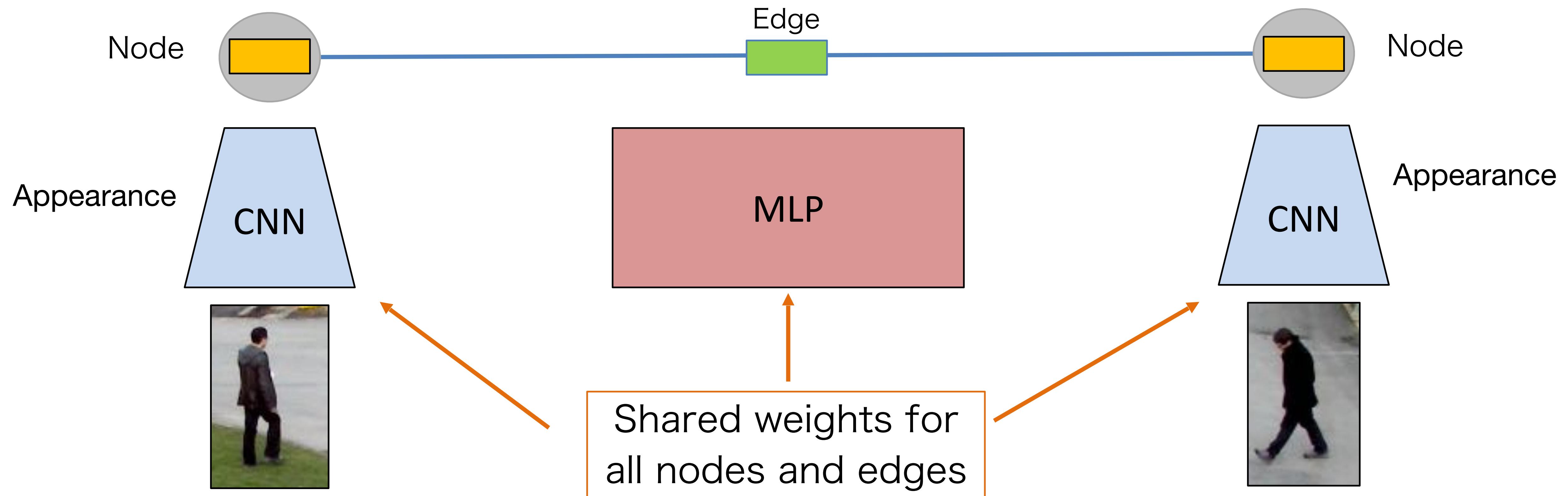
# Feature encoding

- Appearance and geometry encodings



# Feature encoding

- Appearance and geometry encodings



# Feature encoding

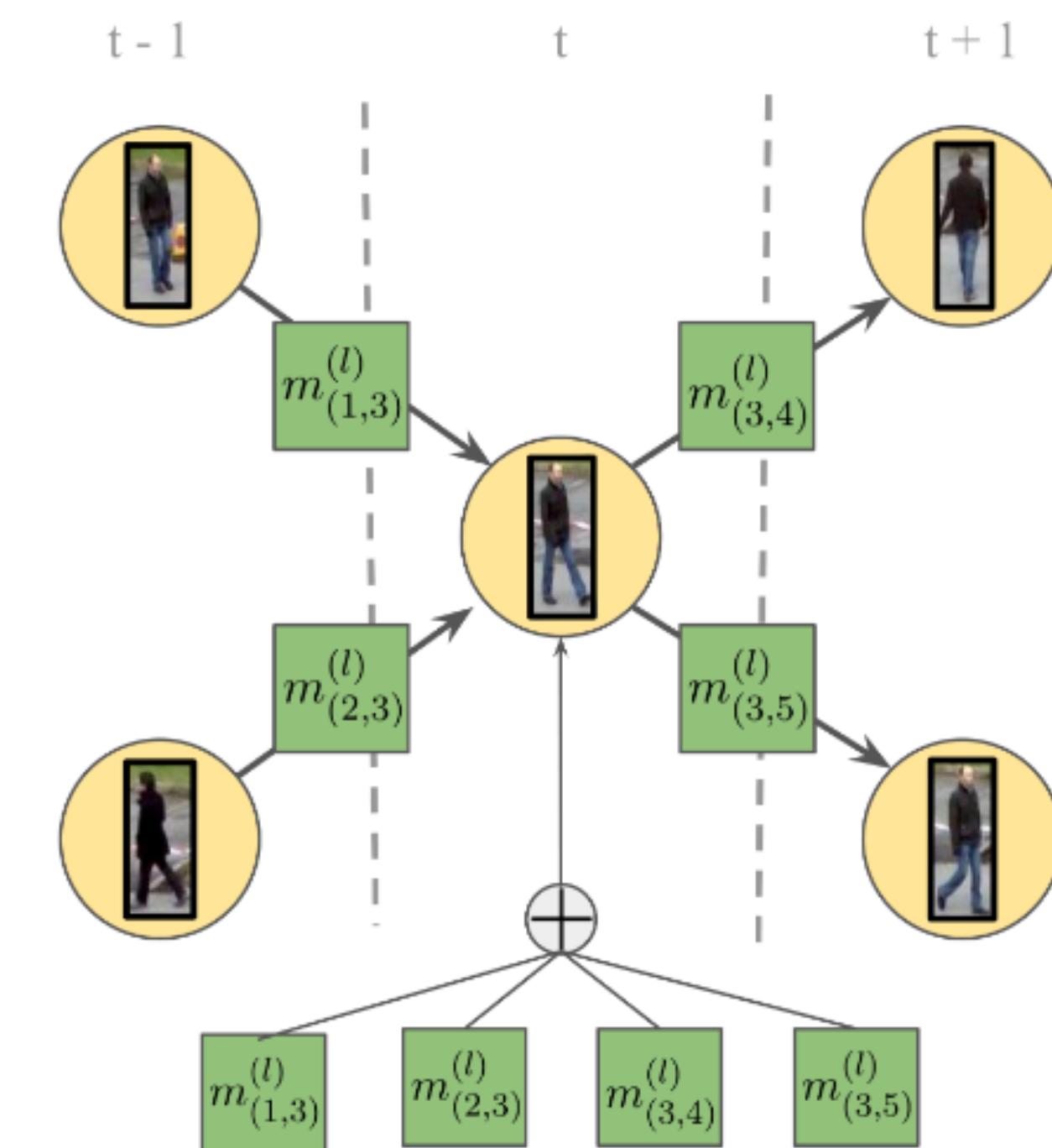
- Contrast:
  - earlier: defining pairwise and unary costs
  - now: feature vectors associated to nodes and edges
- Goal: aggregate context information with multiple iterations of message passing

# Temporal causality

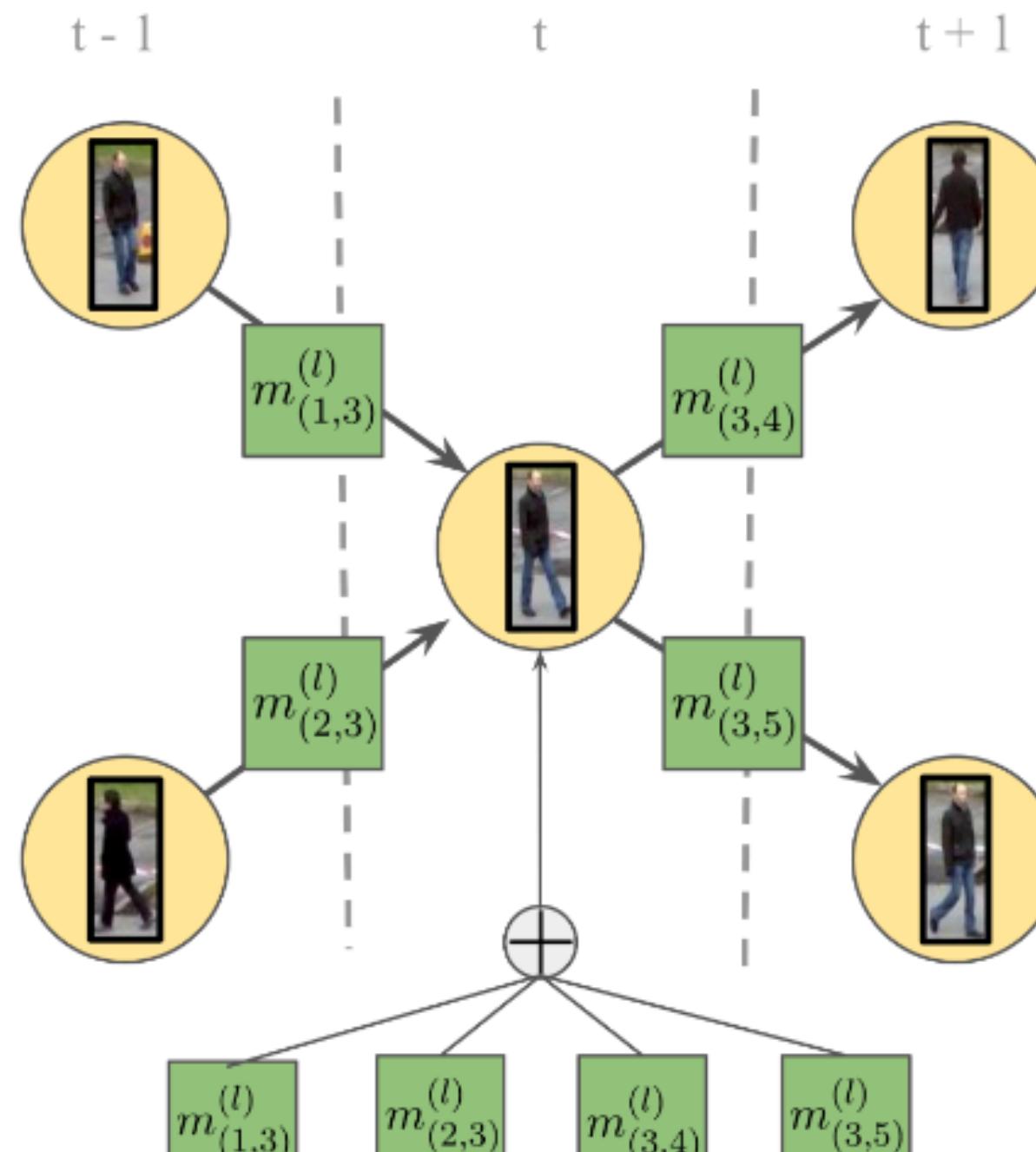
- Recall edge-to-node aggregation:

$$\Phi^{(l)}(i) := \Phi\left(\left\{h^{(l)}(i, j)\right\}_{j \in Ne(i)}\right)$$

- Flow conservation at a node:
  - at most 1 connection to past
  - at most 1 connection to future
- Solution: decouple incident from outgoing edges

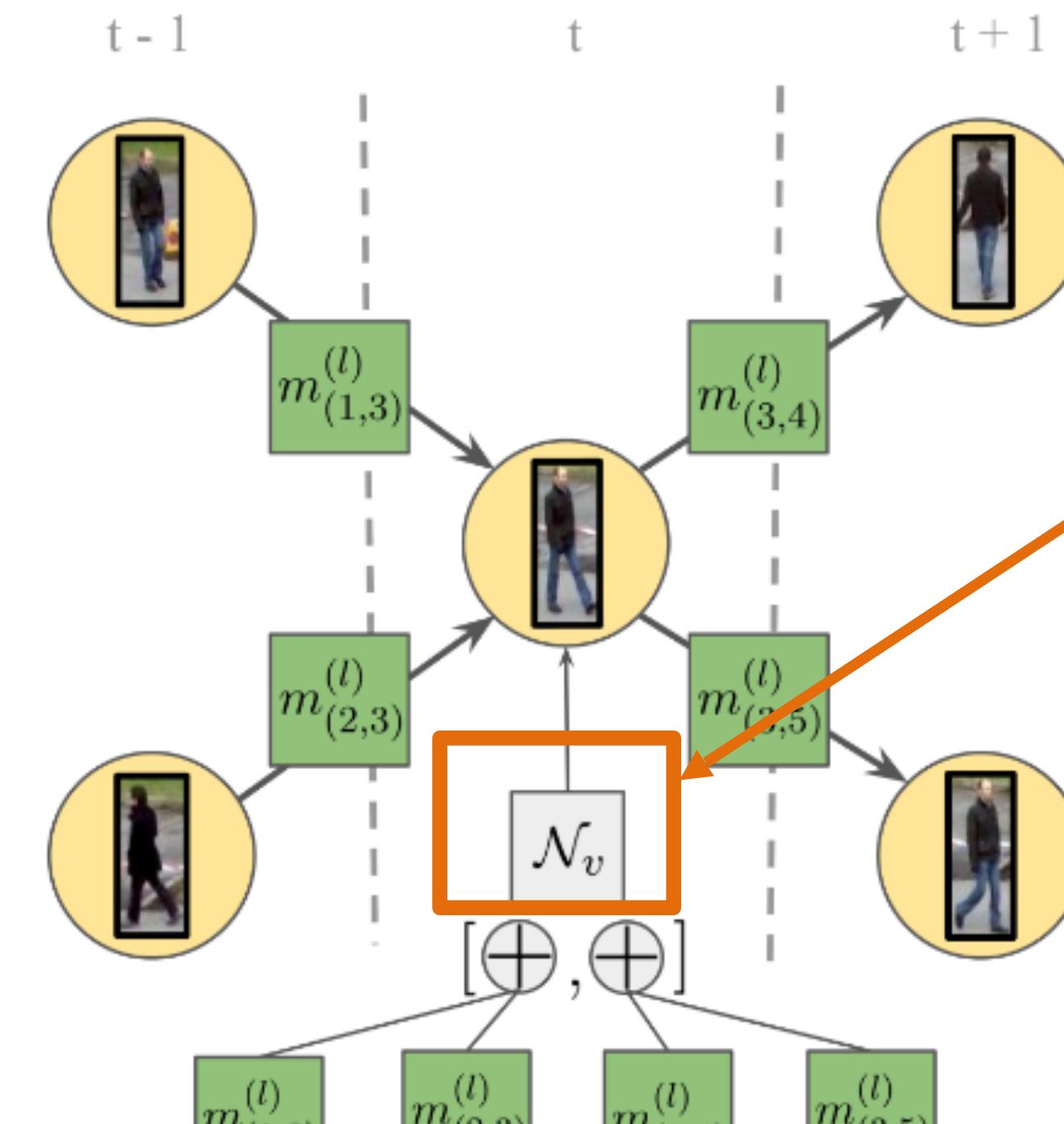


# Time-aware Message Passing



(b) Vanilla node update

All node embeddings are aggregated at once



(c) Time-aware node update

Aggregation of nodes is separated between past / future frames

Brasó and Leal-Taixé. "Learning a Neural Solver for Multiple Object Tracking" (2019).

# Classifying edges

- After several iterations of message passing, each edge embedding contains high-order information about other detections
- We feed the embeddings to an MLP that predicts whether an edge is active/inactive

$$\mathcal{L} = \frac{-1}{|E|} \sum_{l=l_0}^{l=L} \sum_{(i,j) \in E} w \cdot y_{(i,j)} \log(\hat{y}_{(i,j)}^{(l)}) + (1 - y_{(i,j)}) \log(1 - \hat{y}_{(i,j)}^{(l)})$$

Sum over the last steps      Weight to balance active / inactive edges      Edge predictions (w. sigmoid) at iteration l      Binary cross-entropy

# Obtaining final solutions

- After classifying edges, we get a prediction between 0 and 1 for each edge in the graph.
- Directly thresholding solutions could yield infeasible solutions (e.g. violating flow conservation constraints)
- Lightweight post-processing can be used to obtain final trajectories
  - (e.g. linear programming, see Brasó and Leal-Taixé (2019))
- In practice, around 98% of constraints are automatically satisfied, and rounding takes negligible time
- The overall method is fast (~12 fps) and achieves SOTA in the MOT Challenge by a significant margin

# Summary

- No strong assumptions on the graph structure
  - handling occlusions
- Costs can be learned from data
- Accurate and fast (for an offline tracker).
- (Almost) End-to-end learning approach
  - some post-processing required

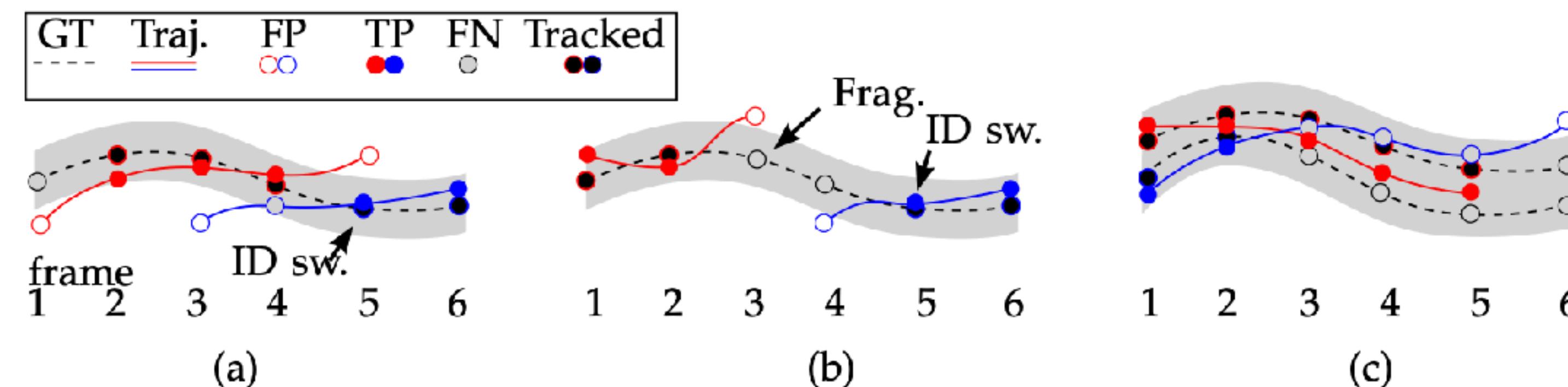
# MOT evaluation

# Evaluation metrics

- Compute a set of measures per frame
  - Perform matching between predictions and ground truth (we will use exactly the same Hungarian algorithm)
  - FP = False positives
  - FN = False negatives (missing detections)
  - IDsw: identity switches

# Evaluation metrics

- How do we compute ID switches?



(a) An ID switch is counted because the ground truth track is assigned first to red, then to blue.

(b) Count both an ID switch (red and blue both assigned to the same ground truth), but also a fragmentation (Frag) because the ground truth coverage was cut.

(c) Identity is preserved. If two trajectories overlap with a ground truth trajectory (within a threshold), the one that forces least ID switches is chosen (the red one).

# Evaluation metrics

- Compute a set of measures per frame
  - Perform matching between predictions and ground truth (we will use exactly the same Hungarian algorithm)
  - FP = False positives
  - FN = False negatives (missing detections)
  - IDsw: identity switches

Multi-object tracking accuracy  $\longrightarrow$  MOTA =  $1 - \frac{\sum_t (\text{FN}_t + \text{FP}_t + \text{IDSW}_t)}{\sum_t \text{GT}_t}$ ,

Ground truth

# Datasets

- MOTChallenge: [www.motchallenge.net](http://www motchallenge net) (people)
  - Several challenges from less to more crowded
- KITTI benchmark: [http://www.cvlibs.net/datasets/kitti/](http://www cvlibs net/datasets/kitti/) (vehicles)
- UA-Detrac: [http://detrac-db.rit.albany.edu](http://detrac-db rit albany edu) (vehicles)





