# EECS 230 Deep Learning

## Lecture 15: Graph Neural Network

Some slides from Simon Prince, Paul-Edouard Sarlin, and Jure Leskovec

# Outline

❑Graph Neural Network

    ❑Graph convolution layer

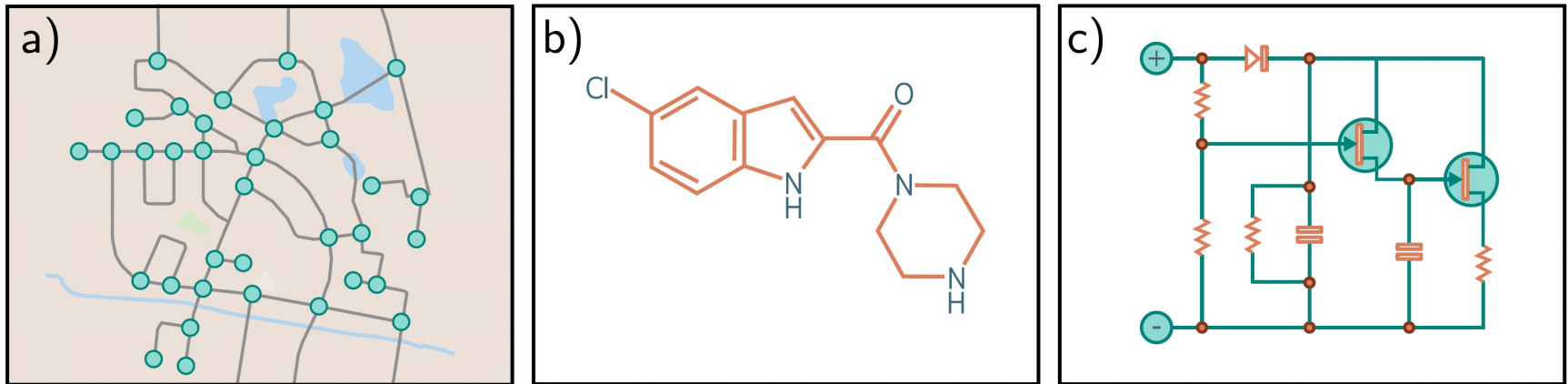    ❑Graph convolutional network

    ❑Graph attention network

❑An application to correspondence matching

    ❑SuperGlue for visual localization

# Graph Neural Network
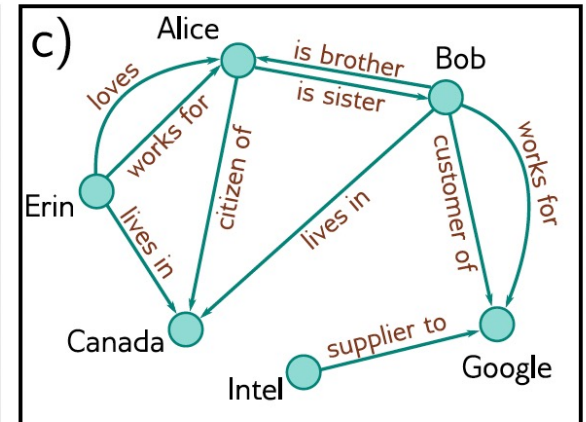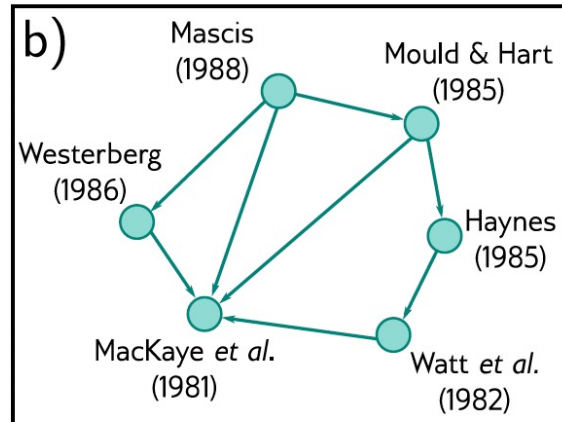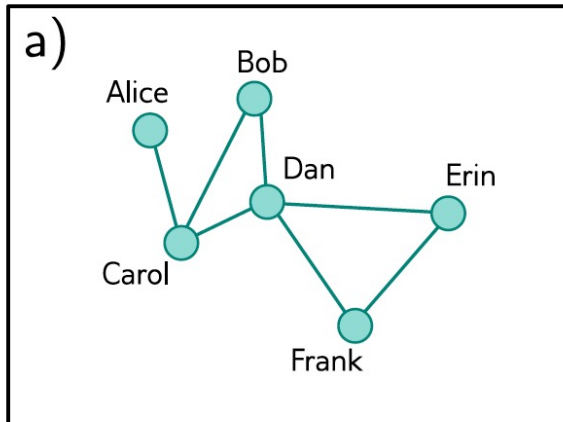
# Real-world graphs



**Figure 13.1** Real-world graphs. Some objects, such as a) road networks, b) molecules, and c) electrical circuits, are naturally structured as graphs.
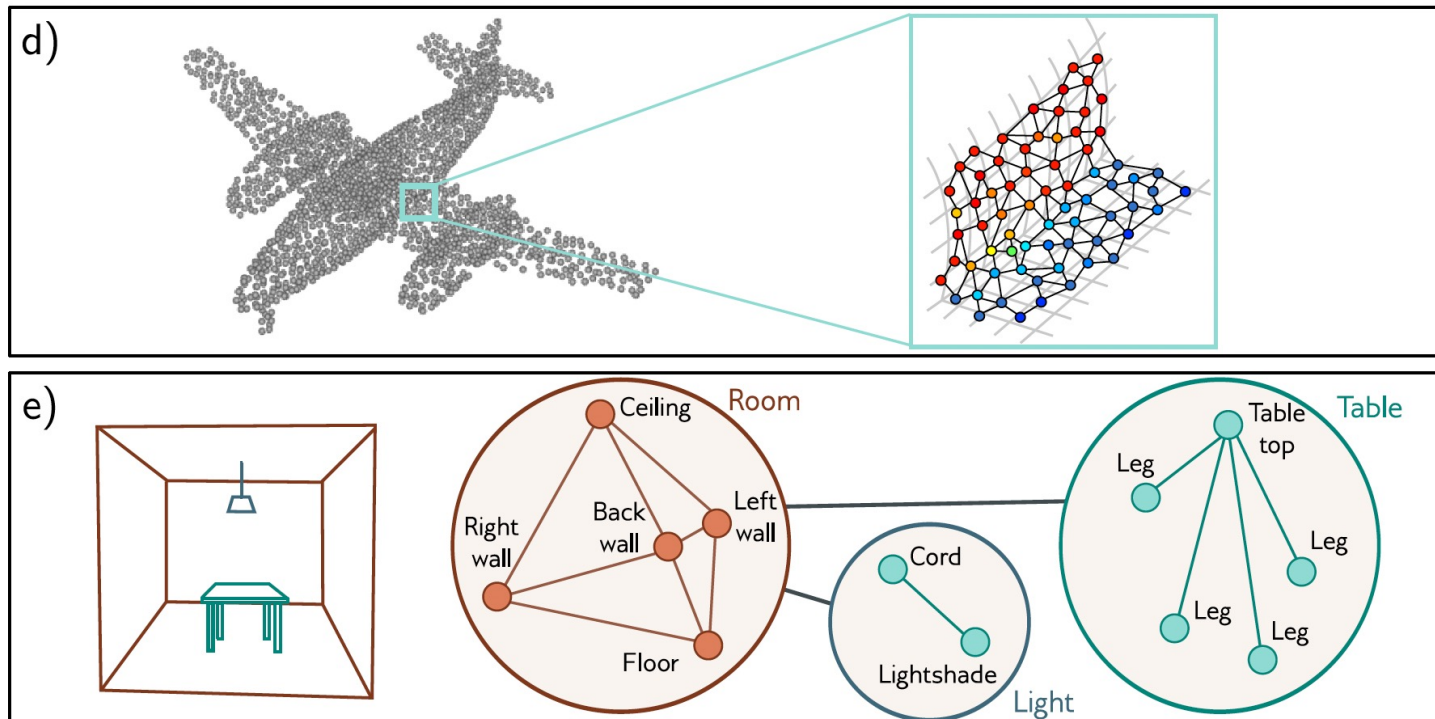
# Types of graphs

❑a) social network is an undirected graph

❑b) citation network is a directed graph

❑c) Knowledge graph is a directed heterogeneous multigraph

# Types of graphs

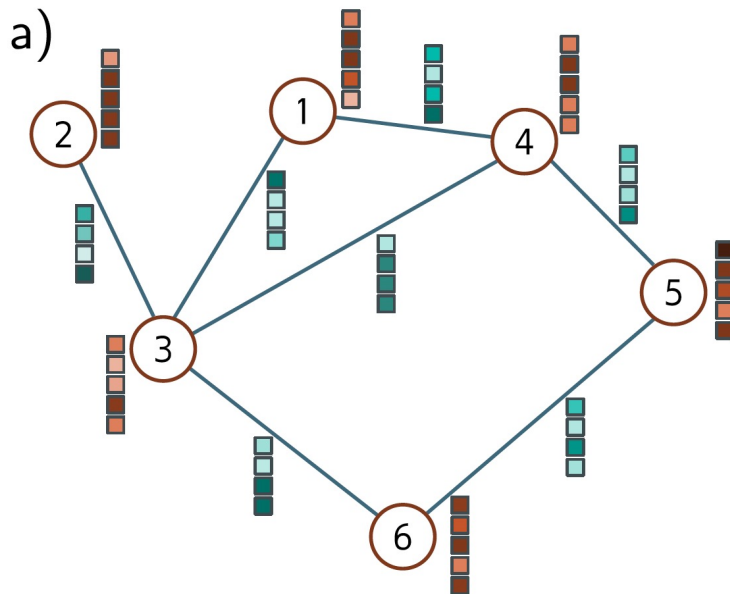❑ d) point cloud as a geometric graph

❑ e) Scene graph is hierarchical

# Graph representation

❑A graph is defined as a tuple G = (V, E)
   ❑where V is a set of nodes
   ❑and E is a set of edges

❑An example graph with 6 nodes and 7 edges

# Properties of adjacency matrix

❑$A^k x$ gives the number of paths of length k to a node

a)



b)

$$A = \begin{bmatrix} 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

c)

$$A^2 = \begin{bmatrix} 2 & 1 & 1 & 1 & 2 & 0 & 0 & 0 \\ 1 & 4 & 1 & 2 & 2 & 1 & 0 & 1 \\ 1 & 1 & 2 & 2 & 1 & 1 & 0 & 1 \\ 1 & 2 & 2 & 3 & 1 & 1 &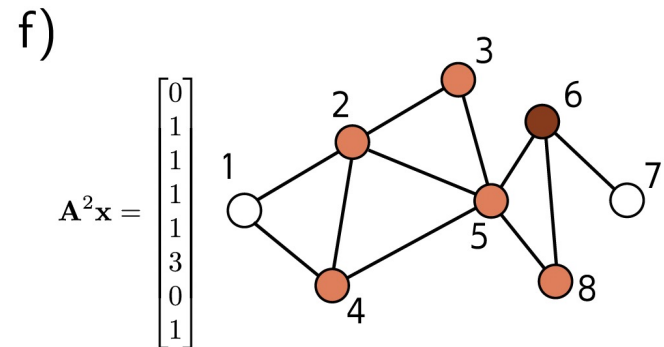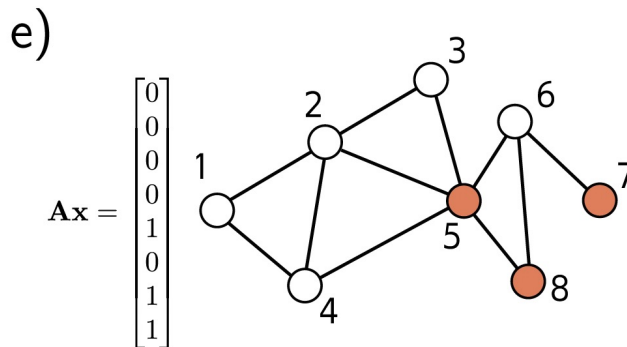 0 & 1 \\ 2 & 2 & 1 & 1 & 5 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 3 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 1 & 1 & 2 \end{bmatrix}$$

d)

$$x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

e)

$$Ax = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$



f)

$$A^2 x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 3 \\ 0 \\ 1 \end{bmatrix}$$

# Permutation invariance

❑ A graph neural network should be permutation invariant

❑ CNN? MLP?

# Tasks on graphs

❑ Graph classification



a)

Graph classification

Graph neural network

Combine    Classify    Class 1
                       Class 2
                       Class 3

# Tasks on graphs

❑Node classification



b)

Node classification

Graph neural network

Classify node from node embedding

# Tasks on graphs

❑ Edge classification



c) Edge prediction

Graph neural network

Predict edge presence from adjacent embeddings

UCMERCED

# Graph convolution

❑ Convolution on a neighborhood



**CNN:** Pixel convolution

**CNN:** Pixel convolution
(as a graph)

**GNN:** Graph convolution

# Graph convolution

Generating features from nodes 1 **hop** away
**A,C,** and **D**

Now embedding for node **B** has information from A,C, and D.

**Processing information from neighbors**

**UCMERCED**

# Graph convolution

Generating features from nodes **1 hop** away
**A,C,** and **D**

**Self loop**

But we don't want to forget
information about B either.

Now embedding for node **B** has
information from A,C,D, **and B
itself**.

**Processing information from neighbors**

# Two layers of graph convolution

Including features from nodes **2 hops** away
**A,B,C,D,E,** and **F**

L=0 represents input node features.

L = 0

L = 0

L = 0

L = 1

L = 0

L = 0

L = 2

L = 1

L = 0

L = 0

L = 0

L = 1

L = 0

L = 0

L = 1

Embedding for node **B** (at layer 2) has information from its first and second hop neighbors.



UCMERCED

# Receptive fields in graph neural networks

❑Increasing receptive fields with more layers/hops

# Graph convolution

Initial 0-th layer embeddings are equal to node features

$$\mathbf{h}_v^0 = \mathbf{x}_v$$

Previous layer embedding of $v$

$$\mathbf{h}_v^k = \sigma\left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1}\right), \quad \forall k \in \{1, ..., K\}$$

$$\mathbf{z}_v = \mathbf{h}_v^K$$

Embedding after K layers of neighborhood aggregation

Non-linearity (e.g., ReLU)

Average of neighbor's previous layer embeddings

UCMERCED

# Graph convolution

- **Recap:** **Simple neighborhood aggregation:**

$$\mathbf{h}_v^k = \sigma \left( \mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right)$$

- Graph convolutional operator:
  - Aggregates messages across neighborhoods, $N(v)$
  - $\alpha_{vu} = 1/|N(v)|$ is the **weighting factor (importance)** of node $u$'s message to node $v$
  - $\Longrightarrow \alpha_{vu}$ is defined **explicitly** based on the structural properties of the graph
  - $\Longrightarrow$ All neighbors $u \in N(v)$ are equally important to node $v$

UCMERCED

# Graph attention network

Can we do better than simple neighborhood aggregation?

Can we let weighting factors $\alpha_{vu}$ to be implicitly defined?

■ **Goal:** Specify arbitrary importances to different neighbors of each node in the graph

■ **Idea:** Compute embedding $h$ of each node in the graph following an **attention strategy:**
  ■ Nodes attend over their neighborhoods' message
  ■ Implicitly specifying different weights to different nodes in a neighborhood

[Velickovic et al., ICLR 2018; Vaswani et al., NIPS 2017]

# Graph attention

- Let $\alpha_{vu}$ be computed as a byproduct of an **attention mechanism $a$:**

  - Let $a$ compute **attention coefficients $e_{vu}$** across pairs of nodes $u$, $v$ based on their messages:
  $$e_{vu} = a(\boldsymbol{W}_k \boldsymbol{h}_u^{k-1}, \boldsymbol{W}_k \boldsymbol{h}_v^{k-1})$$
    - $e_{vu}$ **indicates the importance of node $u'$s message to node $v$**
  - **Normalize coefficients** using the softmax function in order to be comparable across different neighborhoods:
  $$\alpha_{vu} = \frac{\exp(e_{vu})}{\sum_{k \in N(v)} \exp(e_{vk})}$$
  $$\boldsymbol{h}_v^k = \sigma\left(\sum_{u \in N(v)} \alpha_{vu} \boldsymbol{W}_k \boldsymbol{h}_u^{k-1}\right)$$
  **Next:** What is the form of attention mechanism $a$?

UCMERCED

# An application of Graph Neural Network - SuperGlue feature matching

# The importance of context



no
SuperGlue

NN+distance
inliers: 10/29

with
SuperGlue

SuperGlue
inliers: 81/88

# Problem formulation

**Inputs** ➜  ➜ **Outputs**

- Images **A** and **B**
- **2 sets** of **M**, **N** **local features**
  - Keypoints: $\mathbf{p}_i := (x, y, c)_i$
    - Coordinates $(x, y)$
    - Confidence $c$
  - Visual descriptors: $\mathbf{d}_i$

Single a match per keypoint
+ occlusion and noise

$\rightarrow$ a **soft partial assignment**:

$$\mathbf{P} \in [0, 1]^{M \times N}$$

sum ≤ 1

sum ≤ 1

**A Graph Neural Network with attention**

Encodes **contextual cues** & priors

**Reasons** about the 3D scene

**Solving a partial assignment problem**

Differentiable **solver**

Enforces the assignment constraints = **domain knowledge**

**Optimal Matching Layer**

**Attentional Graph Neural Network**

local features

$\mathbf{d}_i^A$ — visual descriptor

$\mathbf{p}_i^A$ — position

$\mathbf{p}_i^B$

$\mathbf{d}_i^B$

Keypoint Encoder

Attentional Aggregation

Self

Cross

matching descriptors
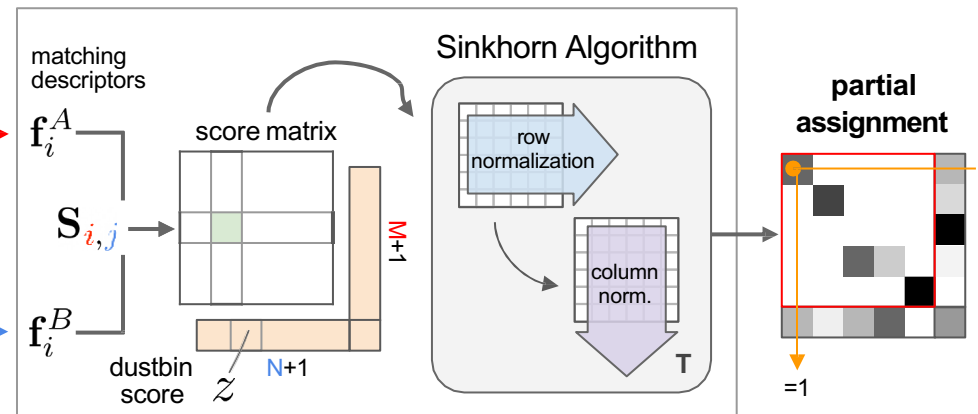
$\mathbf{f}_i^A$

$\mathbf{f}_i^B$

score matrix

$\mathbf{S}_{i,j}$

dustbin score $z$

Sinkhorn Algorithm

row normalization

column norm.

**partial assignment**

=1

- Initial representation for each keypoints $i : {}^{(0)}\mathbf{x}_i$
- Combines visual appearance and position with an MLP:

$$ {}^{(0)}\mathbf{x}_i = \mathbf{d}_i + \mathrm{MLP}\left(\mathbf{p}_i\right) $$

Multi-Layer Perceptron

**Attentional Graph Neural Network**

Attentional Aggregation

local features

$\mathbf{d}_i^A$ visual descriptor

$\mathbf{p}_i^A$ position

$\mathbf{p}_i^B$

$\mathbf{d}_i^B$

Keypoint Encoder

Self    Cross

L

**Optimal Matching Layer**

matching descriptors

$\mathbf{f}_i^A$

$\mathbf{f}_i^B$

Sinkhorn Algorithm

score matrix

$\mathbf{S}_{i,j}$

row normalization

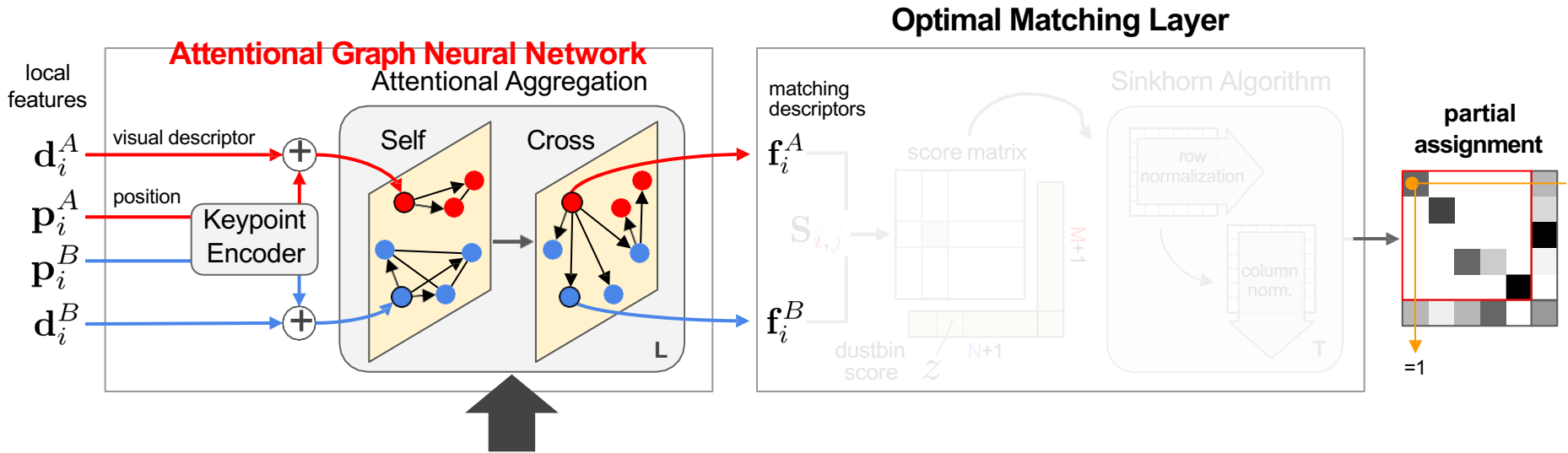column norm.

M+1

N+1

dustbin score $z$

T

**partial assignment**

=1

**Update** the representation based on other keypoints:

- in the same image: "**self**" edges
- in the other image: "**cross**" edges

$\rightarrow$ A complete **graph** with two types of edges

$${}^{(\ell)}\mathbf{x}_i^A \longrightarrow {}^{(\ell+1)}\mathbf{x}_i^A$$

**Optimal Matching Layer**

**Attentional Graph Neural Network**

local features

$\mathbf{d}_i^A$ visual descriptor

$\mathbf{p}_i^A$ position

$\mathbf{p}_i^B$

$\mathbf{d}_i^B$

Keypoint Encoder

Attentional Aggregation

Self   Cross

L

matching descriptors

$\mathbf{f}_i^A$

$\mathbf{f}_i^B$

Sinkhorn Algorithm

score matrix

$\mathbf{S}_{i,j}$

row normalization

column norm.

M+1

dustbin score   $z$   N+1   T

**partial assignment**

=1

**Update** the representation using a **Message Passing Neural Network**

$$^{(\ell+1)}\mathbf{x}_i^A = {}^{(\ell)}\mathbf{x}_i^A + \mathrm{MLP}\left(\left[{}^{(\ell)}\mathbf{x}_i^A \,\|\, \mathbf{m}_{\mathcal{E}\to i}\right]\right)$$

the message

# Attentional Aggregation

- Compute the **message** $\mathbf{m}_{\mathcal{E} \to i}$ using **self** and **cross attention**
- Soft database retrieval: query $\mathbf{q}_i$, key $\mathbf{k}_j$, and value $\mathbf{v}_j$

$$\mathbf{m}_{\mathcal{E} \to i} = \sum_{j:(i,j) \in \mathcal{E}} \alpha_{ij} \mathbf{v}_j \qquad \mathbf{q}_i = \mathbf{W}_1 {}^{(\ell)} \mathbf{x}_i + \mathbf{b}_1$$

$$\alpha_{ij} = \text{Softmax}_j \left( \mathbf{q}_i^\top \mathbf{k}_j \right) \qquad \begin{bmatrix} \mathbf{k}_j \\ \mathbf{v}_j \end{bmatrix} = \begin{bmatrix} \mathbf{W}_2 \\ \mathbf{W}_3 \end{bmatrix} {}^{(\ell)} \mathbf{x}_j + \begin{bmatrix} \mathbf{b}_2 \\ \mathbf{b}_3 \end{bmatrix}$$
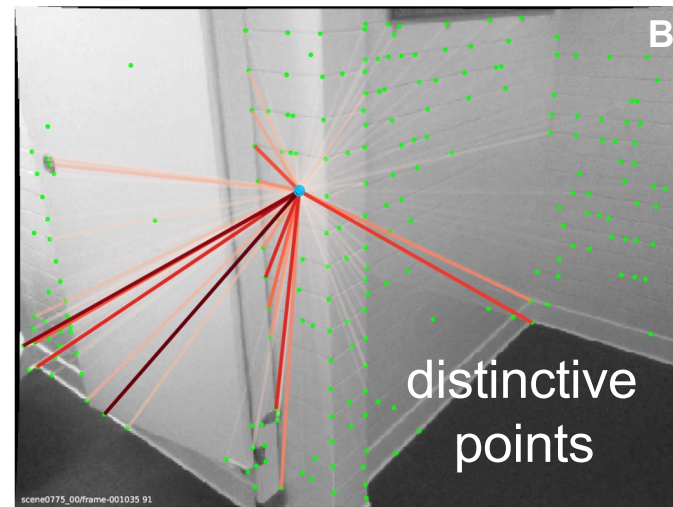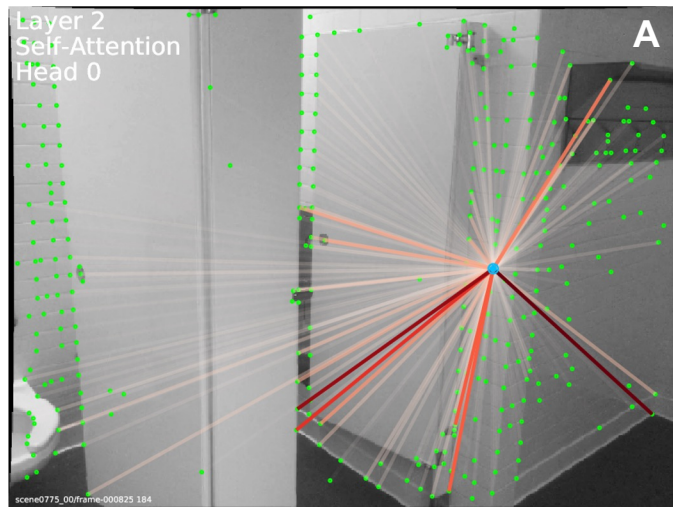
query neighbors

query salient points

= [tile, pos. (80, 110)]

= [corner, pos. (60, 90)]

= [grid, pos. (400, 600)]
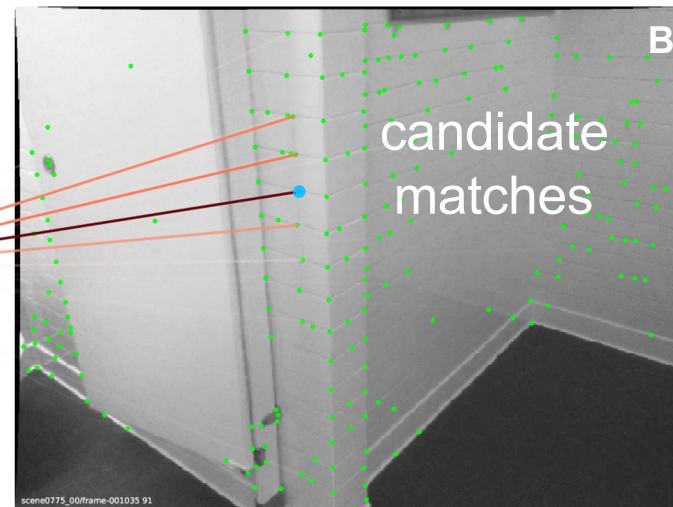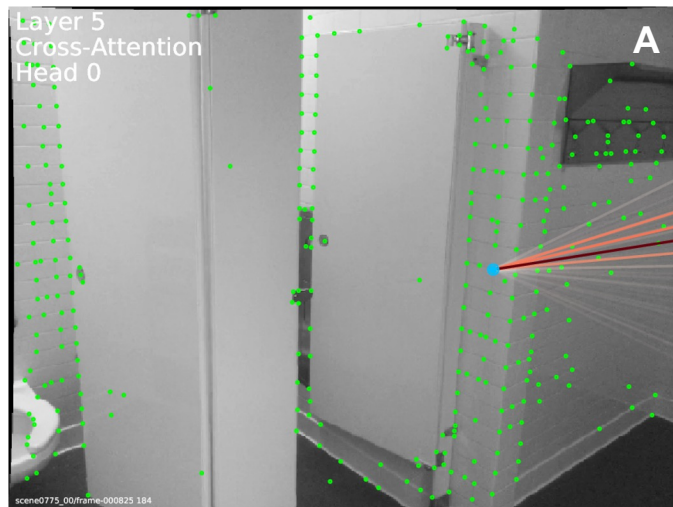
$\mathbf{x}_i$ = [tile, position (70, 100)]
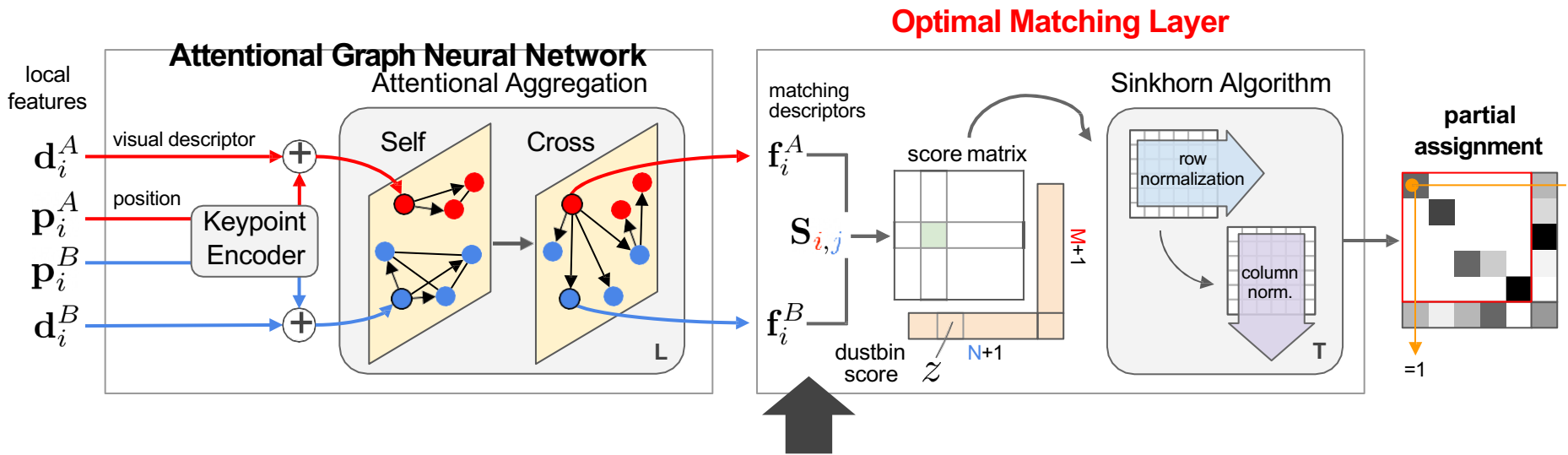
[Vaswani et al, 2017]

**Self-attention**
= intra-image
information
flow
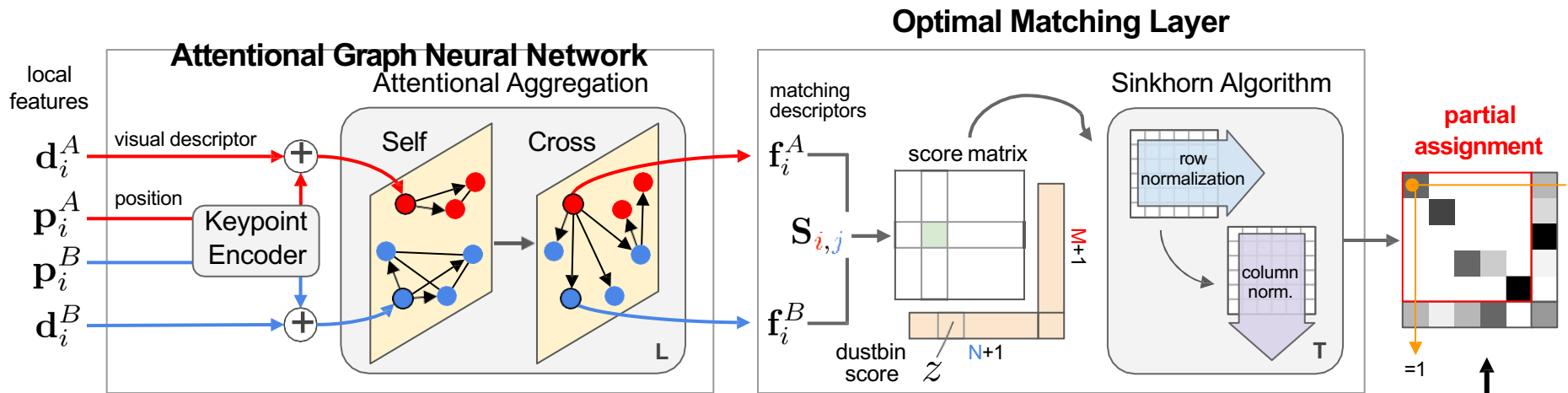
**Cross-attention**
= inter-image

Attention builds a
**soft**, **dynamic**,
**sparse graph**

Compute a **score matrix** $\mathbf{S} \in \mathbb{R}^{M \times N}$ for all matches:

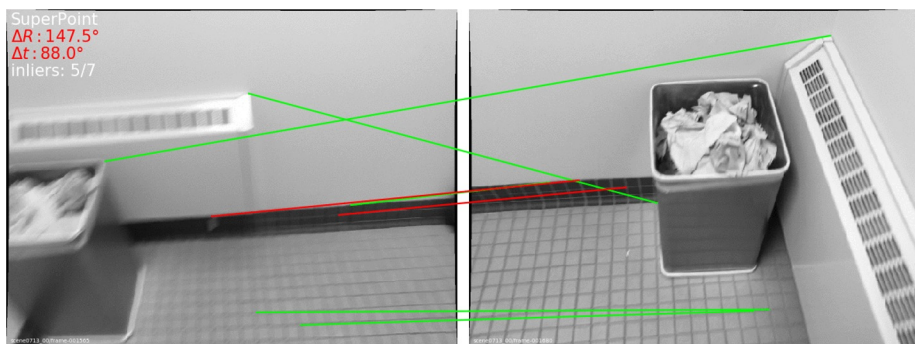$$\mathbf{f}_i^A = \mathbf{W} \cdot {}^{(L)}\mathbf{x}_i^A + \mathbf{b}$$
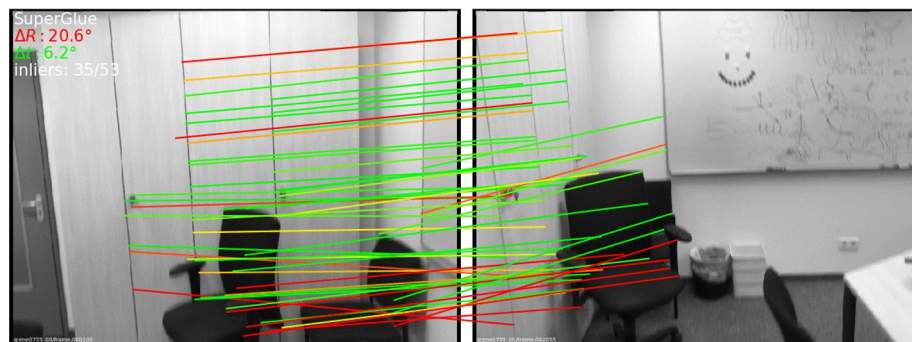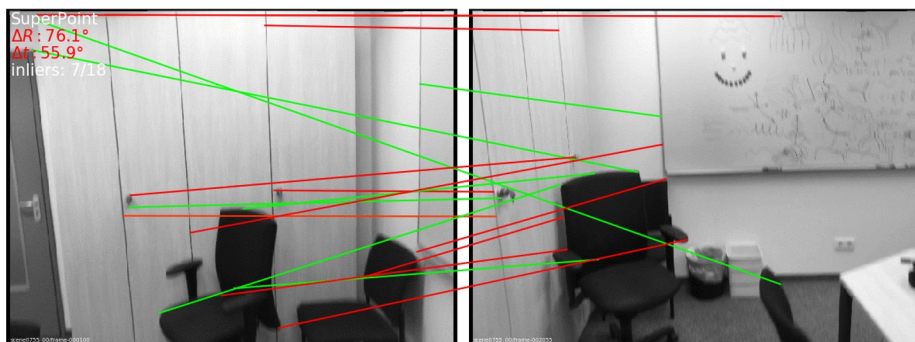
$$\mathbf{S}_{i,j} = <\mathbf{f}_i^A, \mathbf{f}_j^B>$$

**Attentional Graph Neural Network**

Attentional Aggregation

local features

$\mathbf{d}_i^A$   visual descriptor

$\mathbf{p}_i^A$   position

$\mathbf{p}_i^B$

$\mathbf{d}_i^B$

Keypoint Encoder

Self   Cross

L

**Optimal Matching Layer**

matching descriptors

$\mathbf{f}_i^A$

$\mathbf{S}_{i,j}$

$\mathbf{f}_i^B$

score matrix

M+1

dustbin score   $z$   N+1

Sinkhorn Algorithm

row normalization

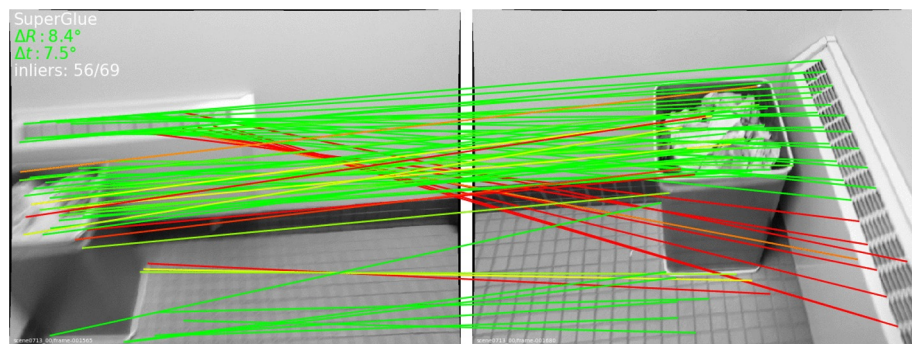column norm.

T

**partial assignment**

=1

- Compute **ground truth correspondences** from pose and depth
- Find which keypoints should be **unmatched**
- Loss: maximize the log-likelihood $\bar{\mathbf{P}}_{i,j}$ of the GT cells

# Results: indoor - ScanNet

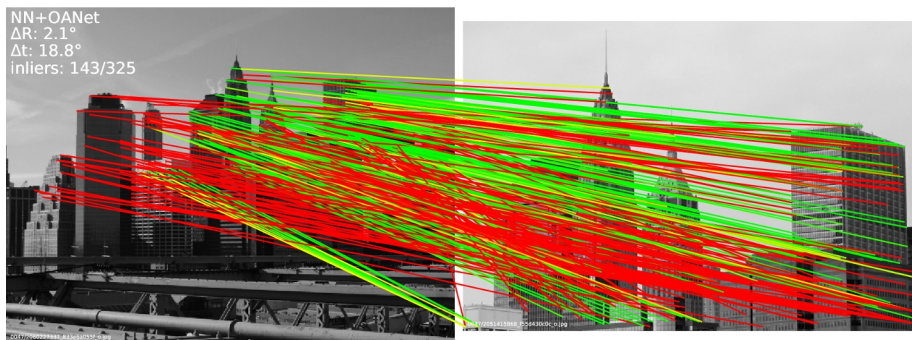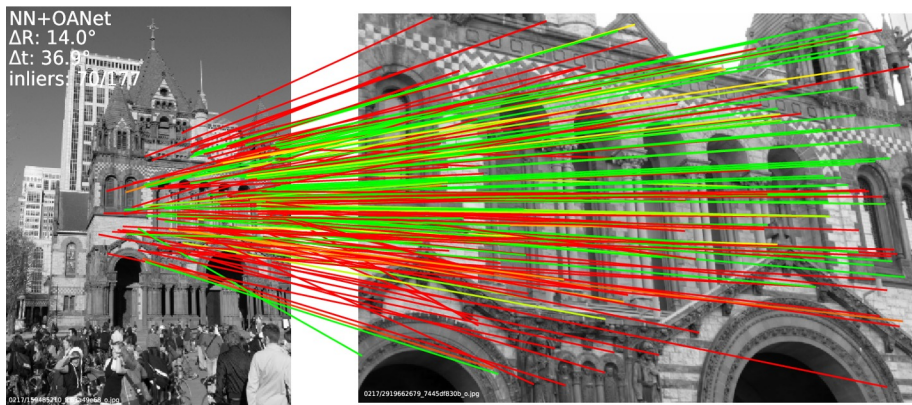SuperPoint + NN + heuristics

SuperPoint + **SuperGlue**



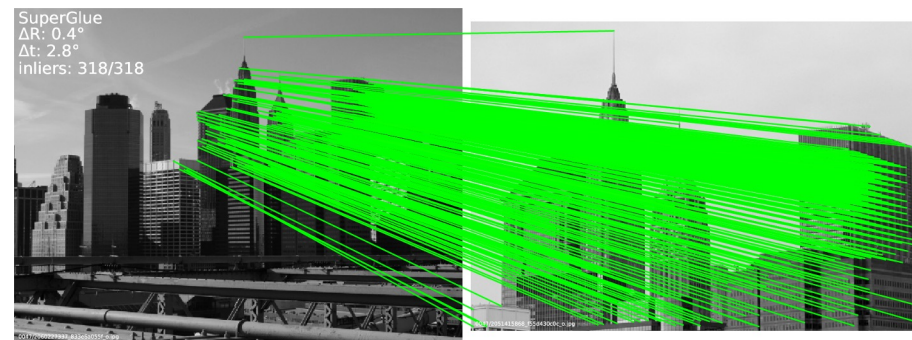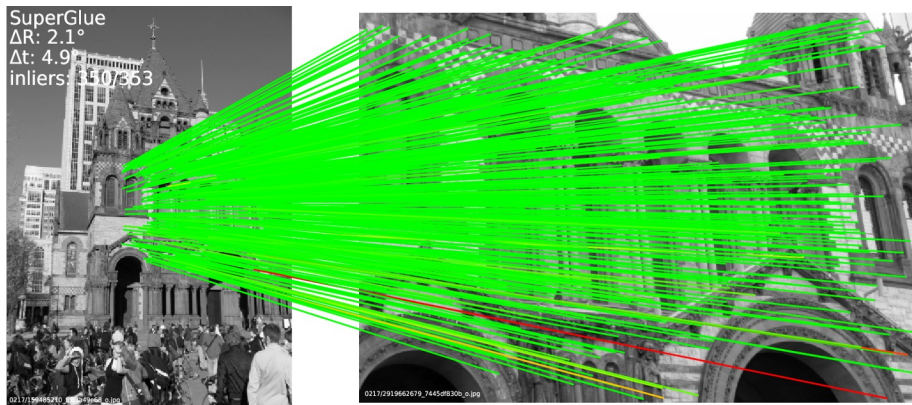SuperGlue: more **correct matches** and fewer **mismatches**

# Results: outdoor - SfM

SuperPoint + NN + OA-Net (inlier classifier)

SuperPoint + **SuperGlue**



SuperGlue: more **correct matches** and fewer **mismatches**