

---

# Robust Learning-based Image Matching

IMW CVPR 2019

---

Dawei Sun

Zixin Luo

Jiahui Zhang

# 00 Outline

An image matching pipeline: 1) local keypoint detection, 2) local keypoint description, 3) sparse matching.



# Outline

An image matching pipeline: 1) local keypoint detection, 2) local keypoint description, 3) sparse matching.

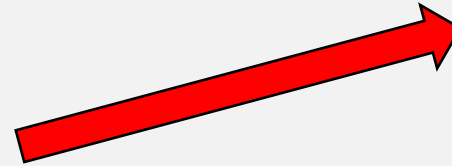
Part 1

**ContextDesc: a learning-based  
local descriptor**

ContextDesc: Local Descriptor

Augmentation with Cross-Modality

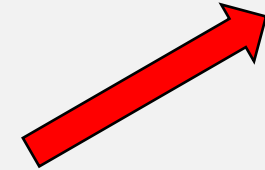
Context, CVPR'19



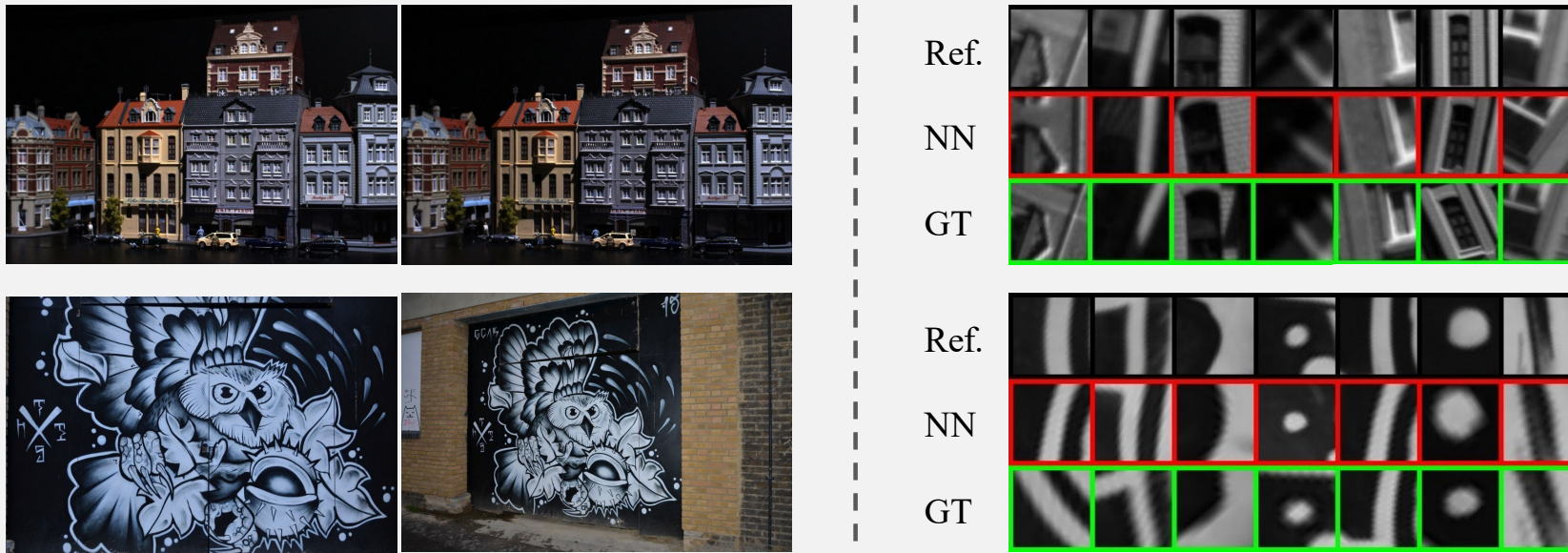
Part 2

**A learning-based inlier  
classification and fundamental  
matrix estimation method**

In submission



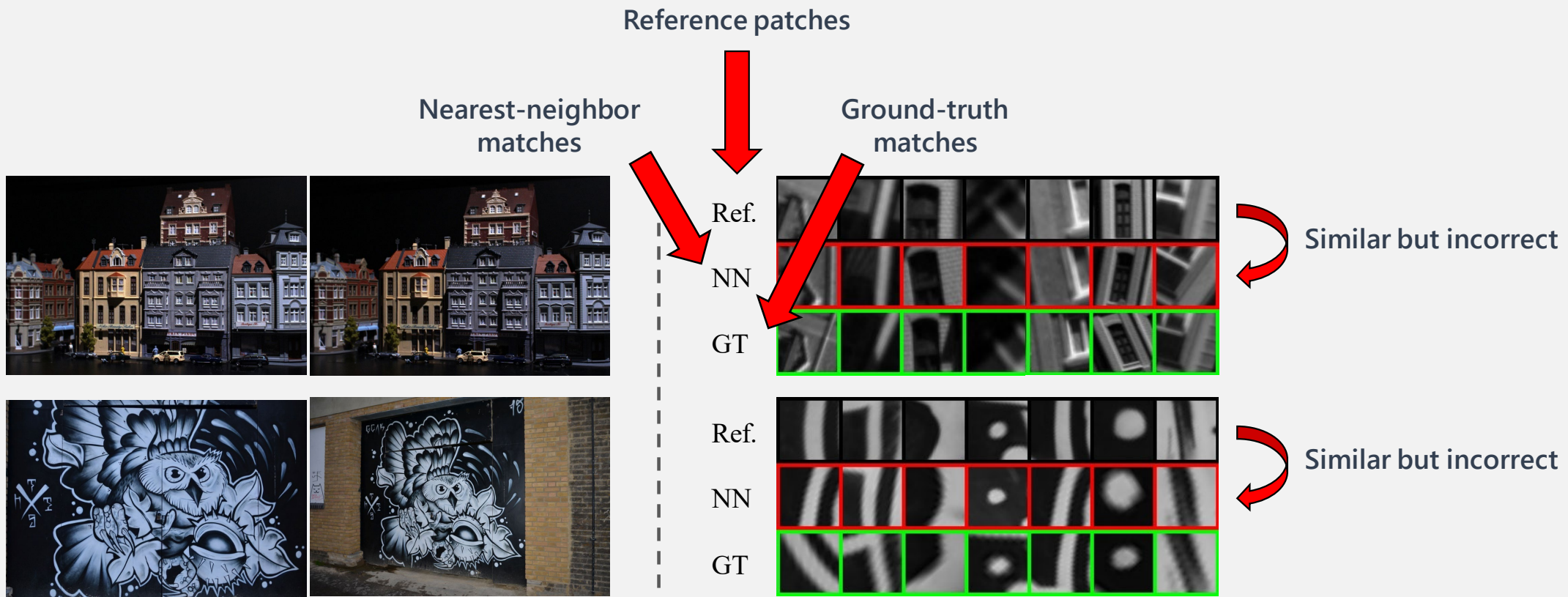
# 01 Motivation



The results are becoming saturated on standard benchmark

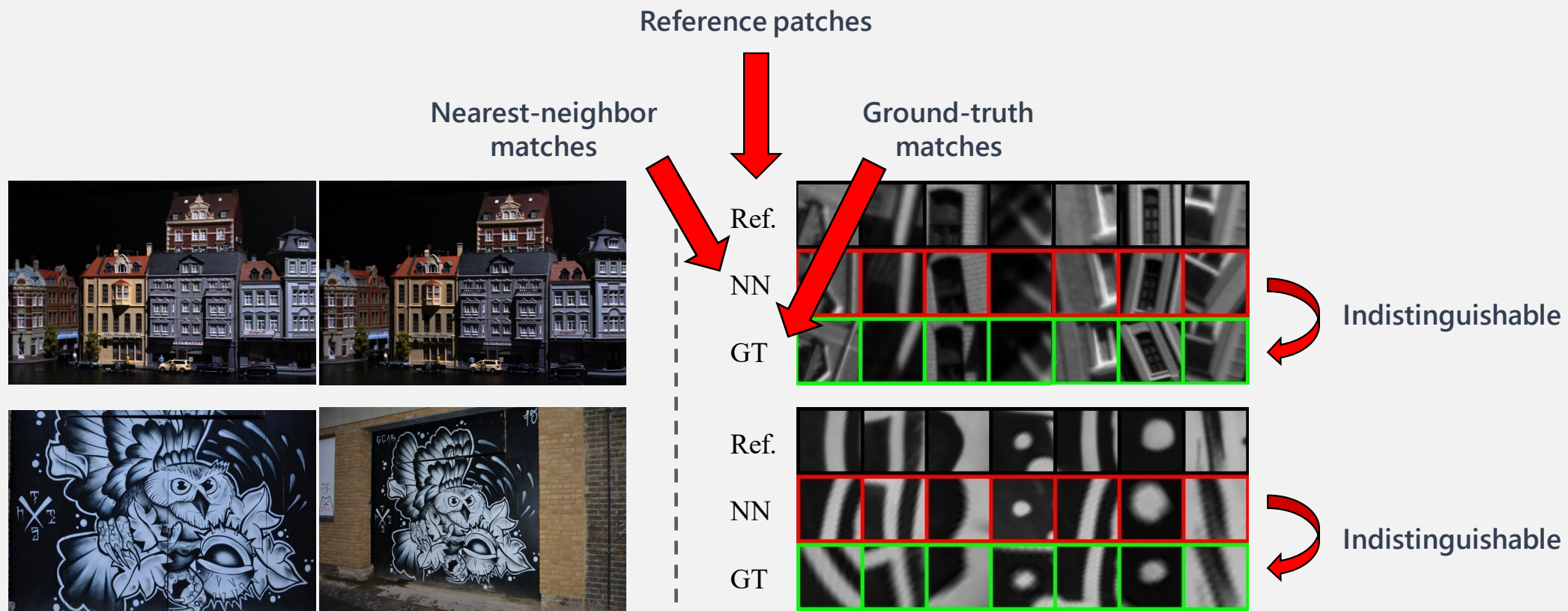
\*Samples are from HPatches dataset.

# 01 Motivation



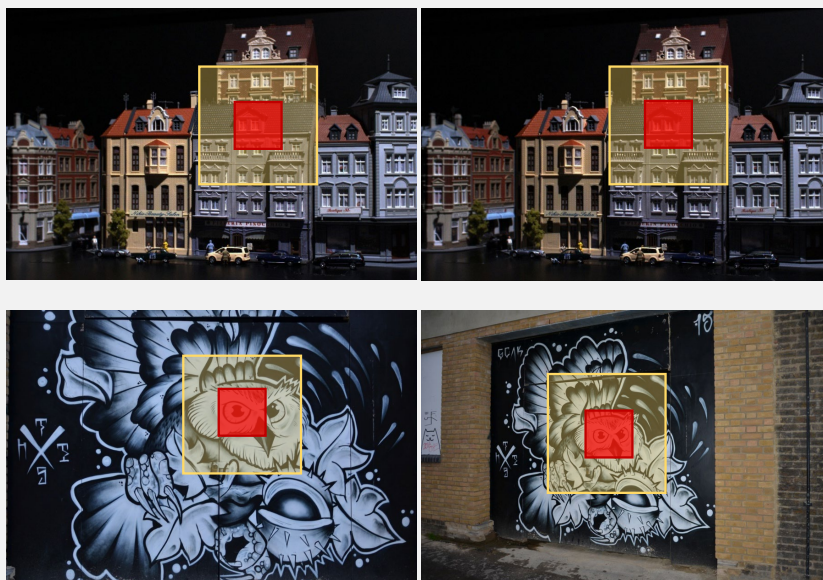
\*GeoDesc is used for feature description.

# 01 Motivation



Locally distinguishable by visual appearance for,  
e.g., repetitive patterns?

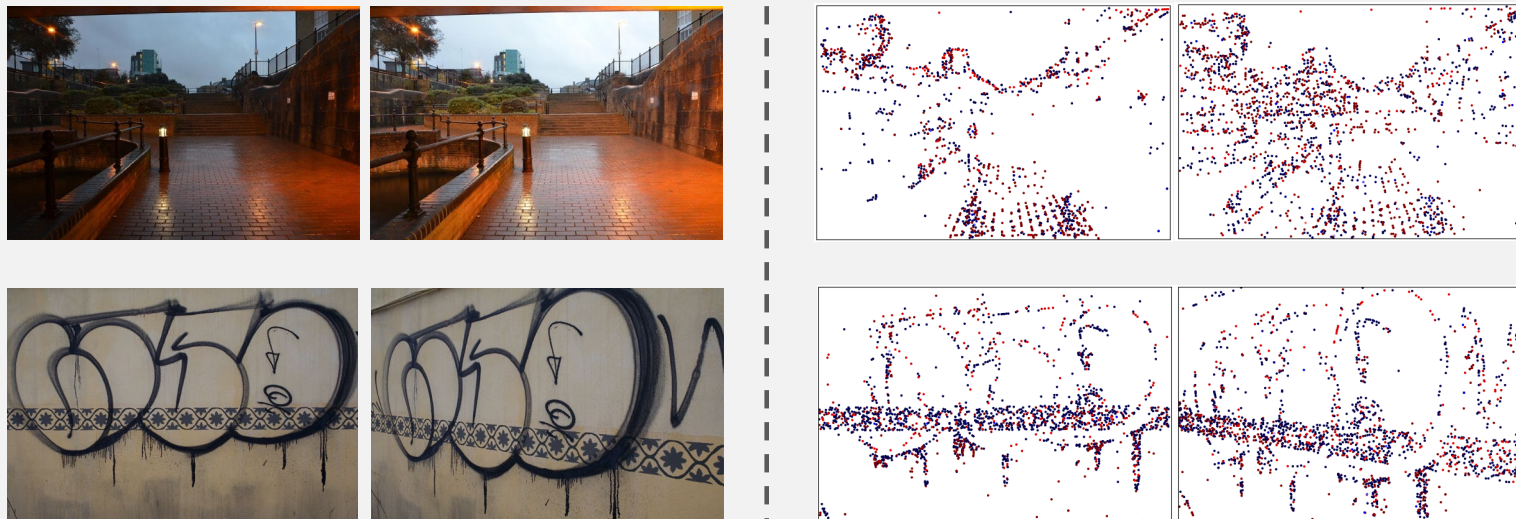
# 01 Motivation



More *visual context!*

- The representation of both local details and richer context – *how to construct the network?*
- Construct feature pyramids – *too costly for this low-level task?*

# 01 Motivation

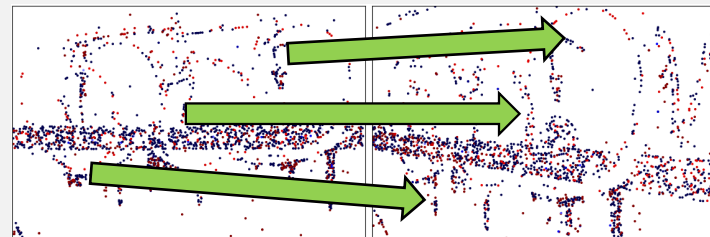
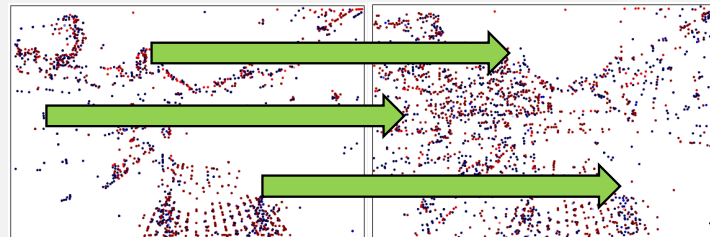
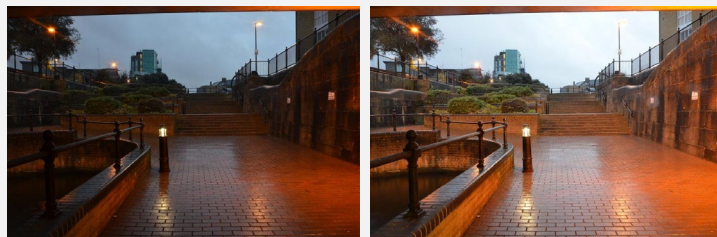


Keypoint distribution reveals meaningful scene structure

\*Keypoints are derived from SIFT.



# 01 Motivation

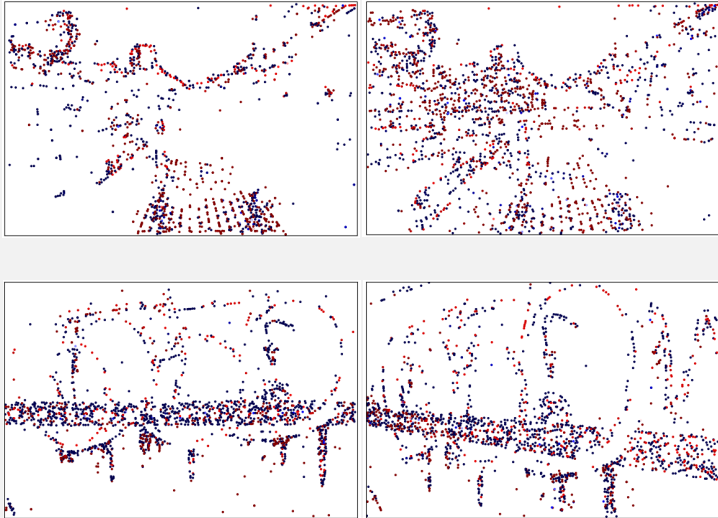


*Coarse* matches can be established, even *without* color information

Keypoint distribution reveals meaningful scene structure

*Keypoints are designed to be repeatable in the same underlying scene*

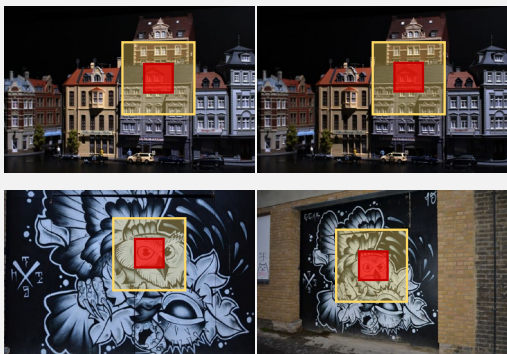
# 01 Motivation



Encoding *geometric context* from keypoint distribution of individual image

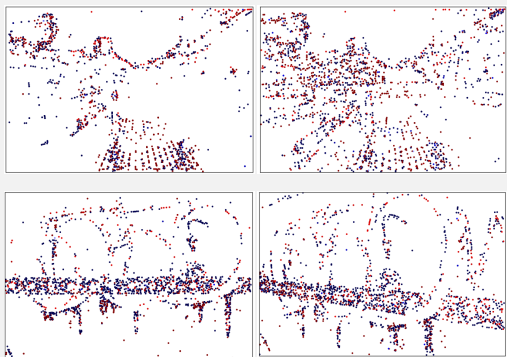
- Keypoints are irregular and unordered – *how to construct a proper encoder?*
- Keypoints are not perfectly repeatable – *how to acquire strong invariance property to different image variations?*

# 01 Motivation



## Visual context

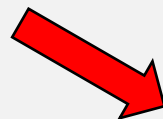
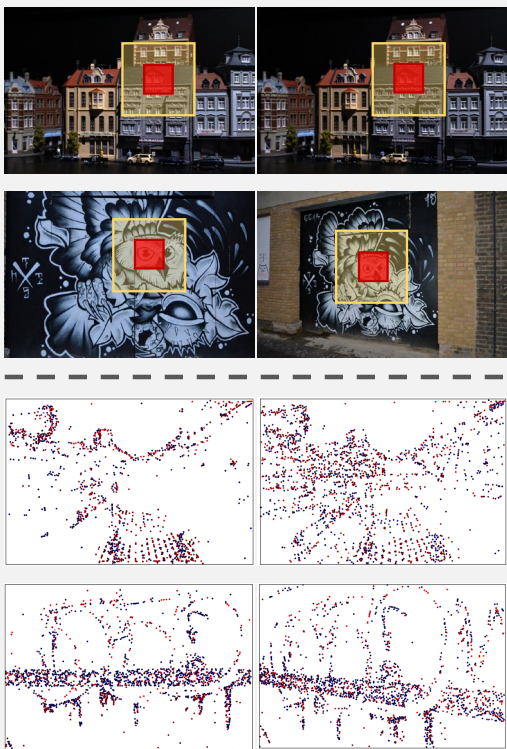
- Incorporate high-level visual information
- Resort to *regional representation* often used in image retrieval (one forward pass of the entire image)



## Geometric context

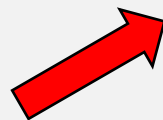
- Geometric cues from keypoint distribution.
- Resort to *PointNet-like architecture* to process 2D point sets

# 01 Motivation

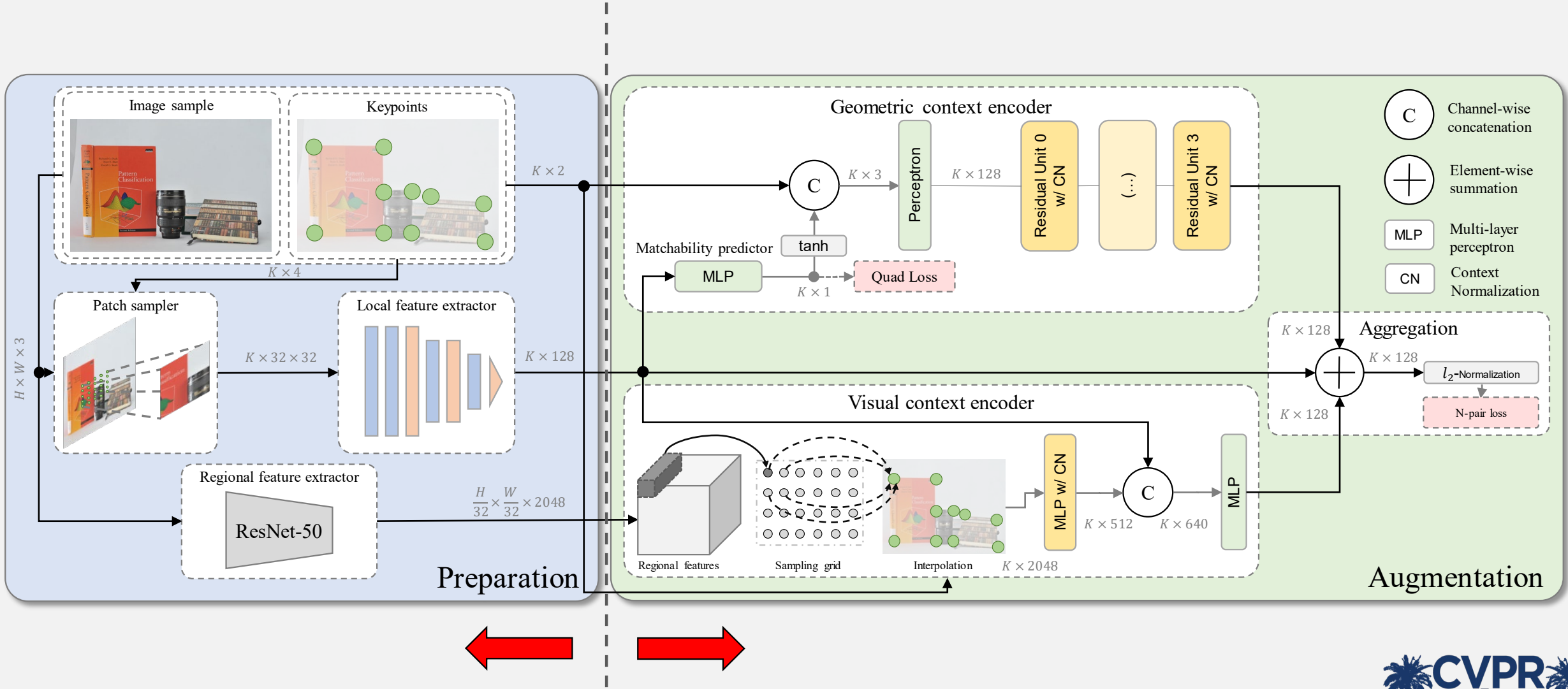


*Based on off-the-shelf descriptors...*

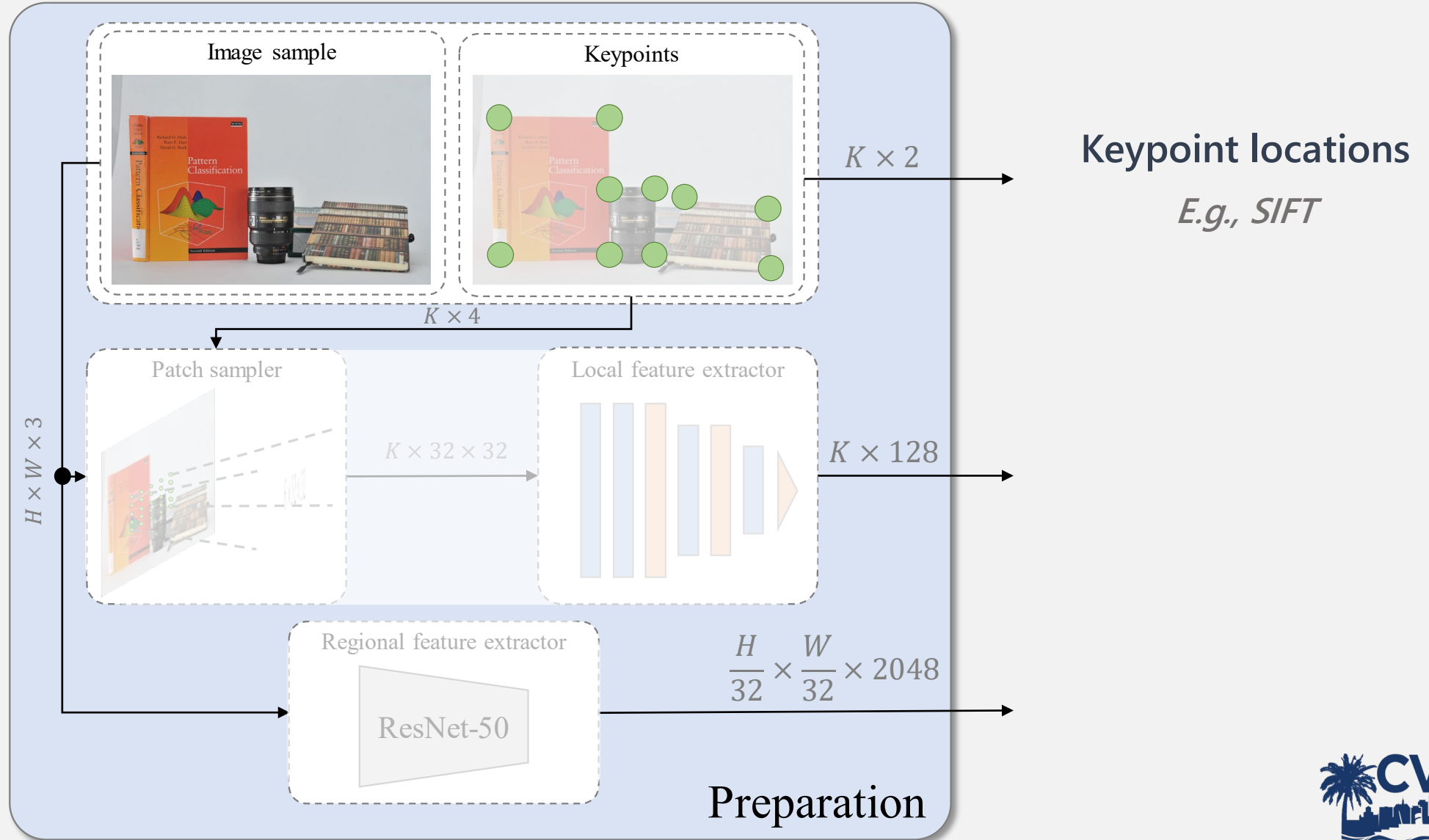
**A unified framework: Cross-modality local descriptor augmentation**



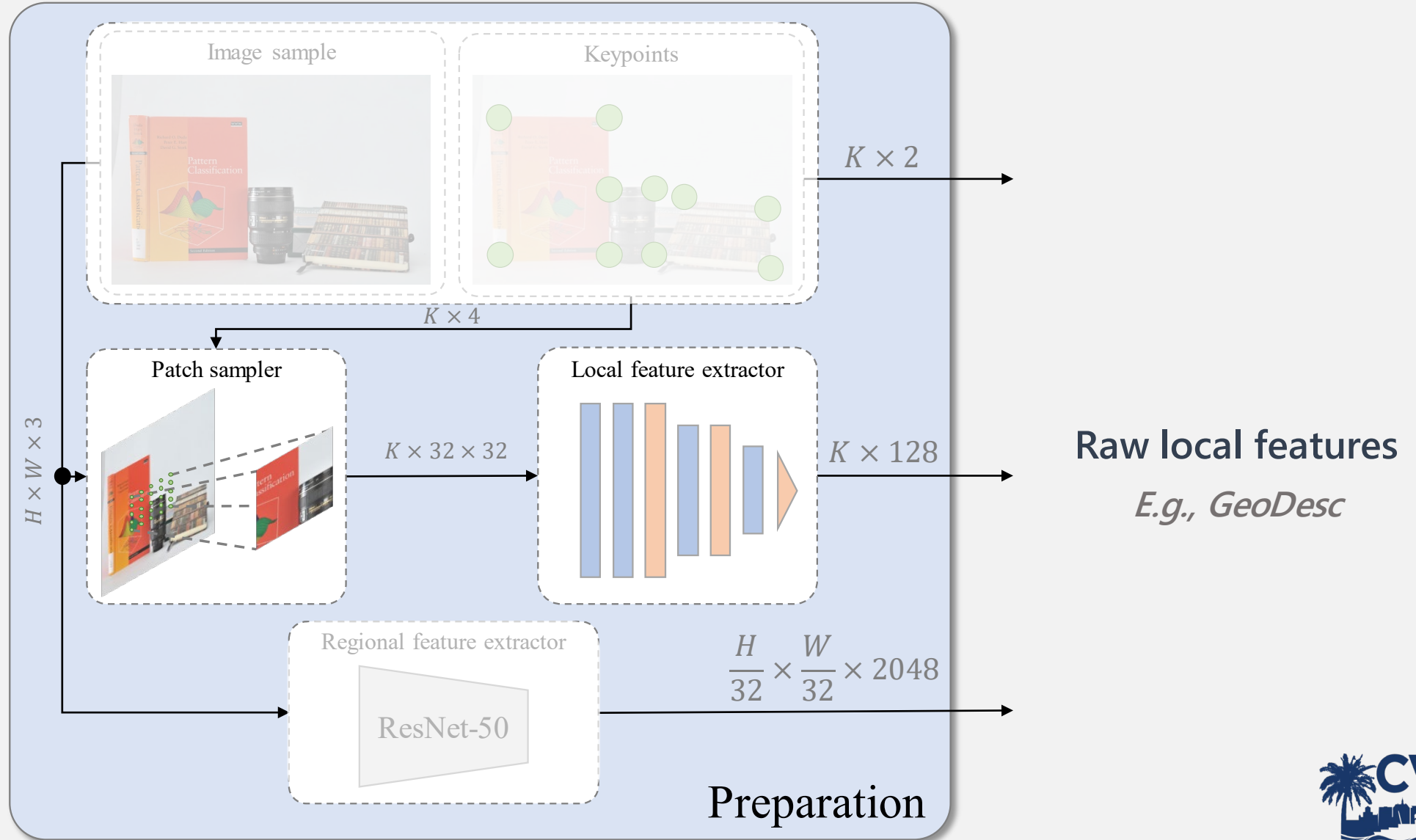
# 02 Methods



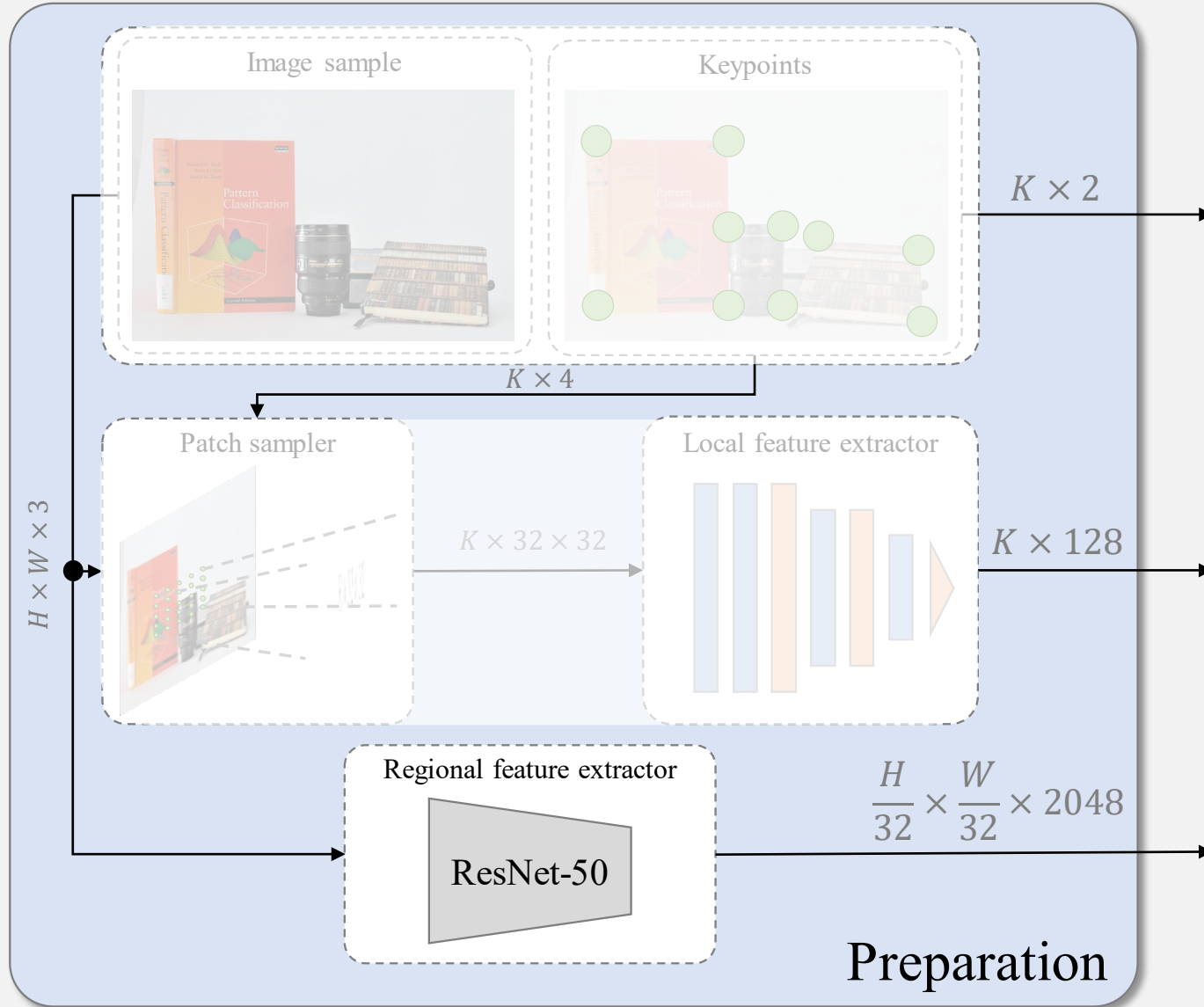
# 02 Methods



# 02 Methods



# 02 Methods



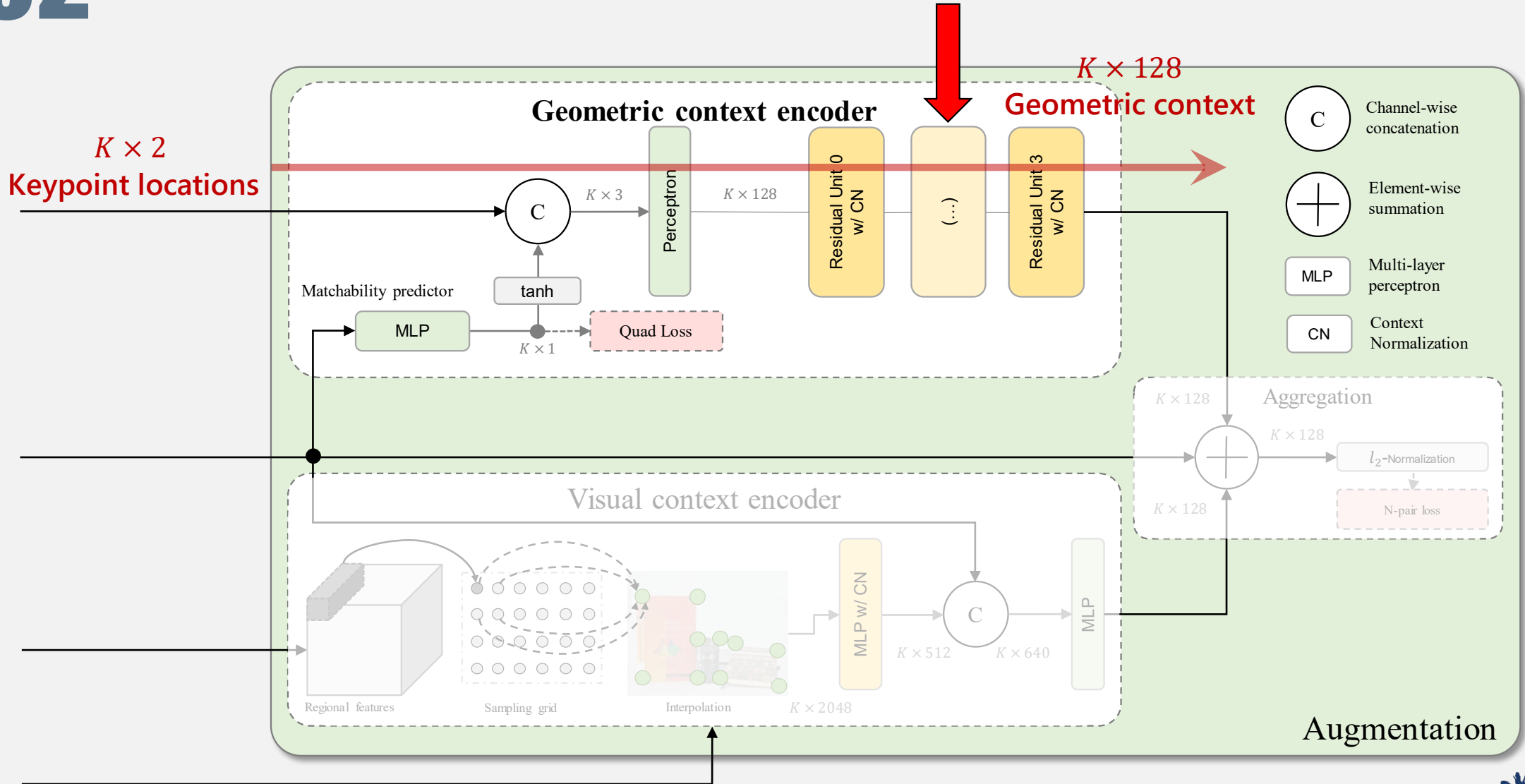
**Regional features**  
*An off-the-shelf image  
retrieval model*



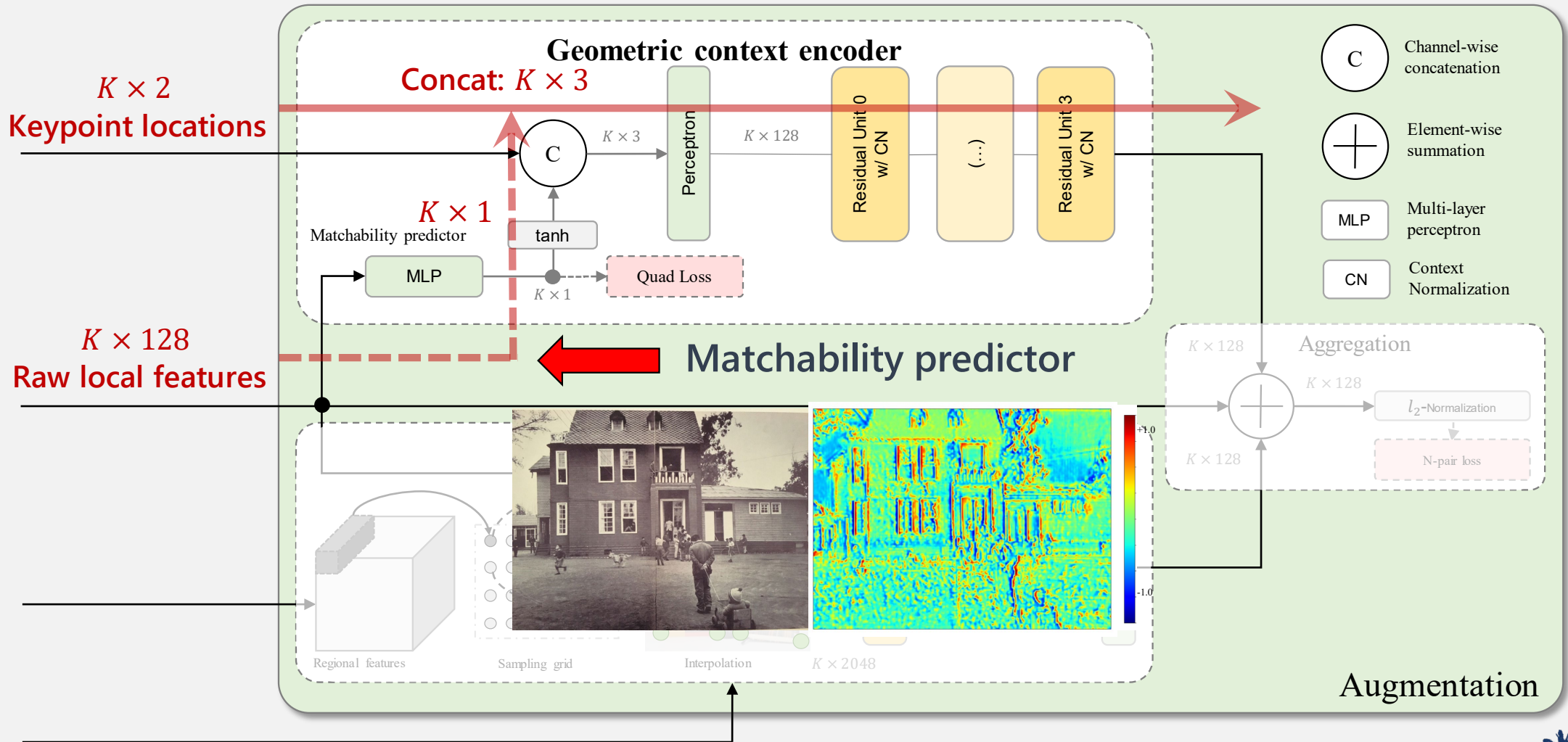


# 02 Methods

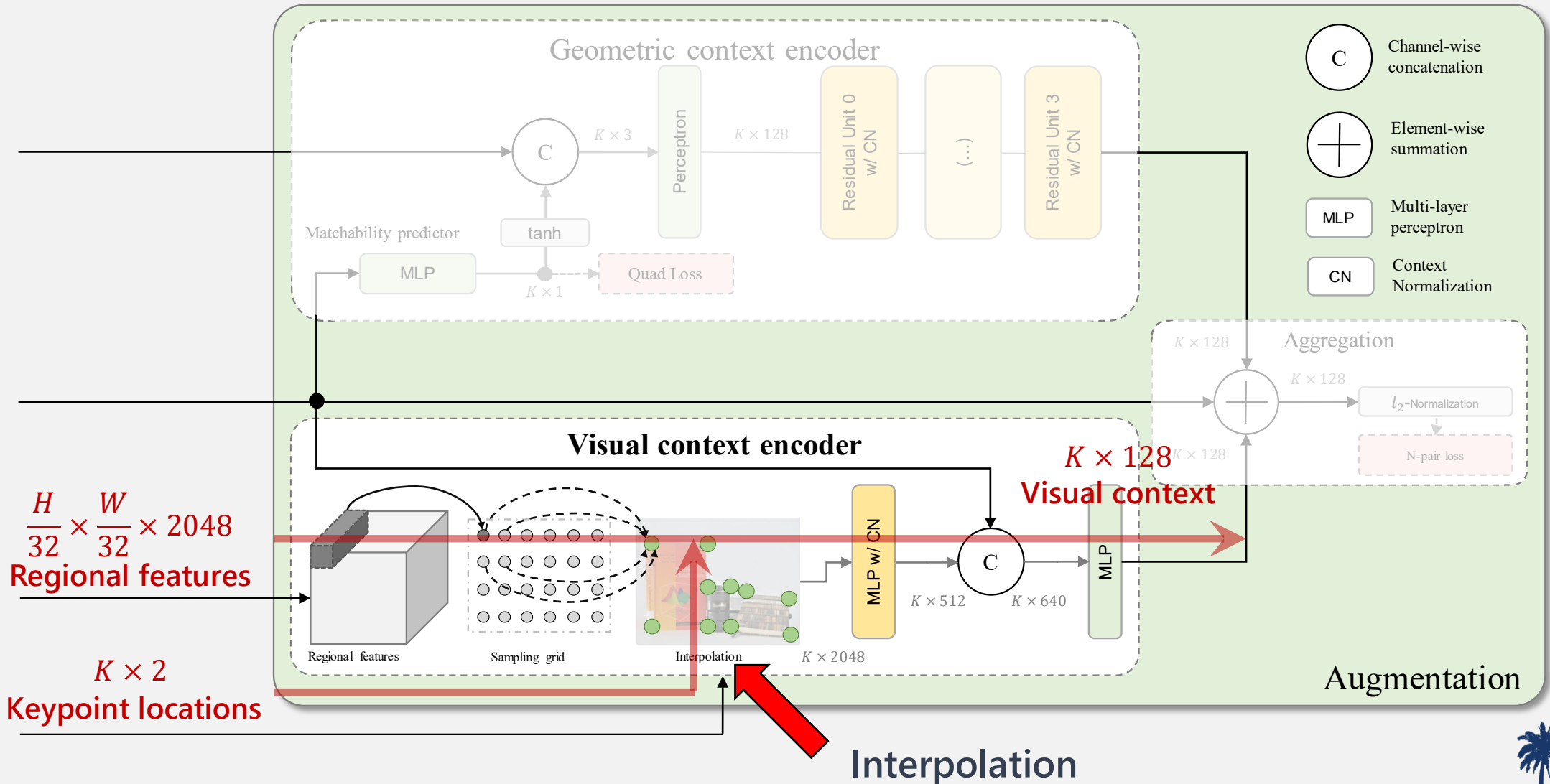
## A variant of PointNet



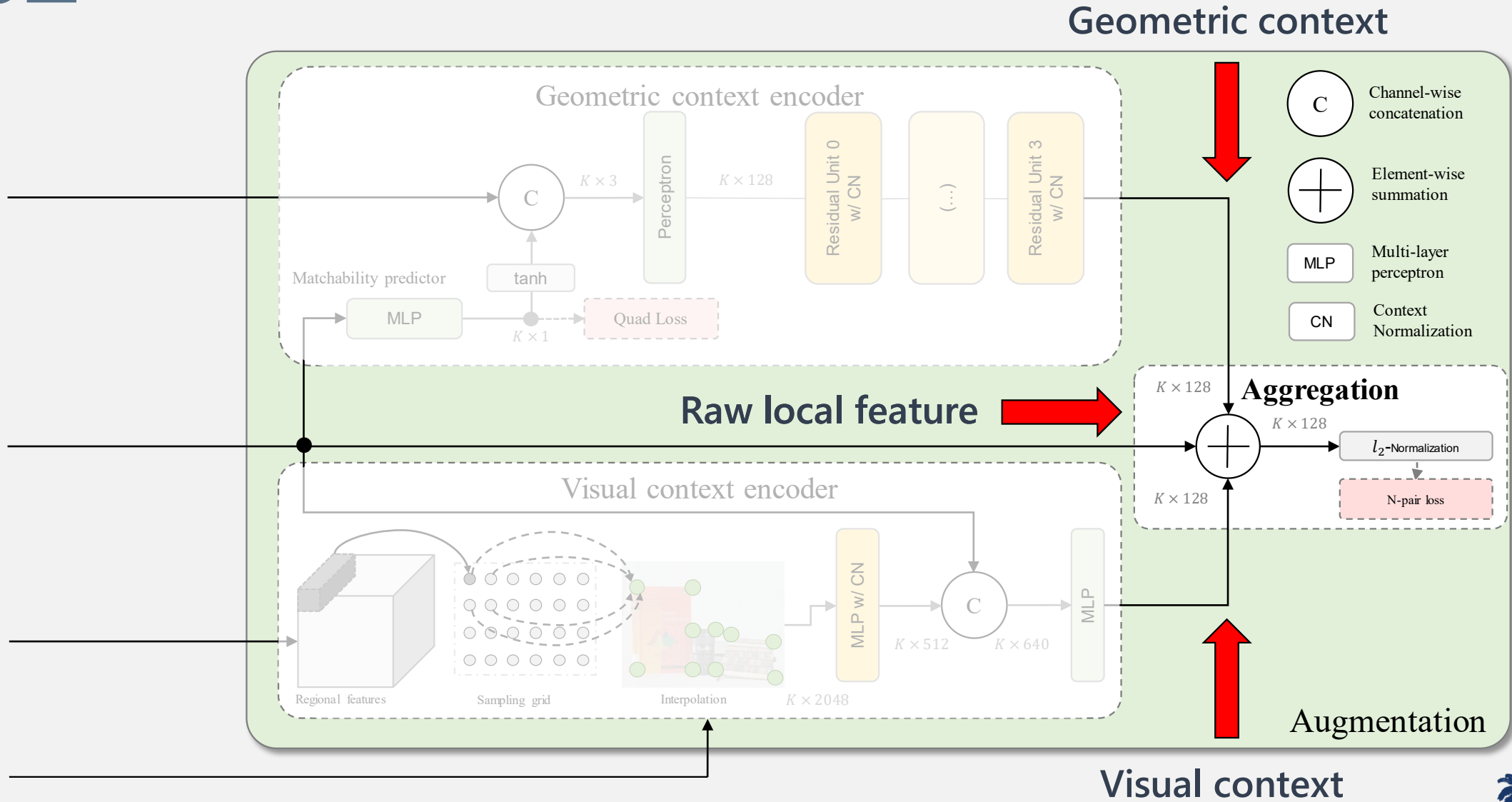
# 02 Methods



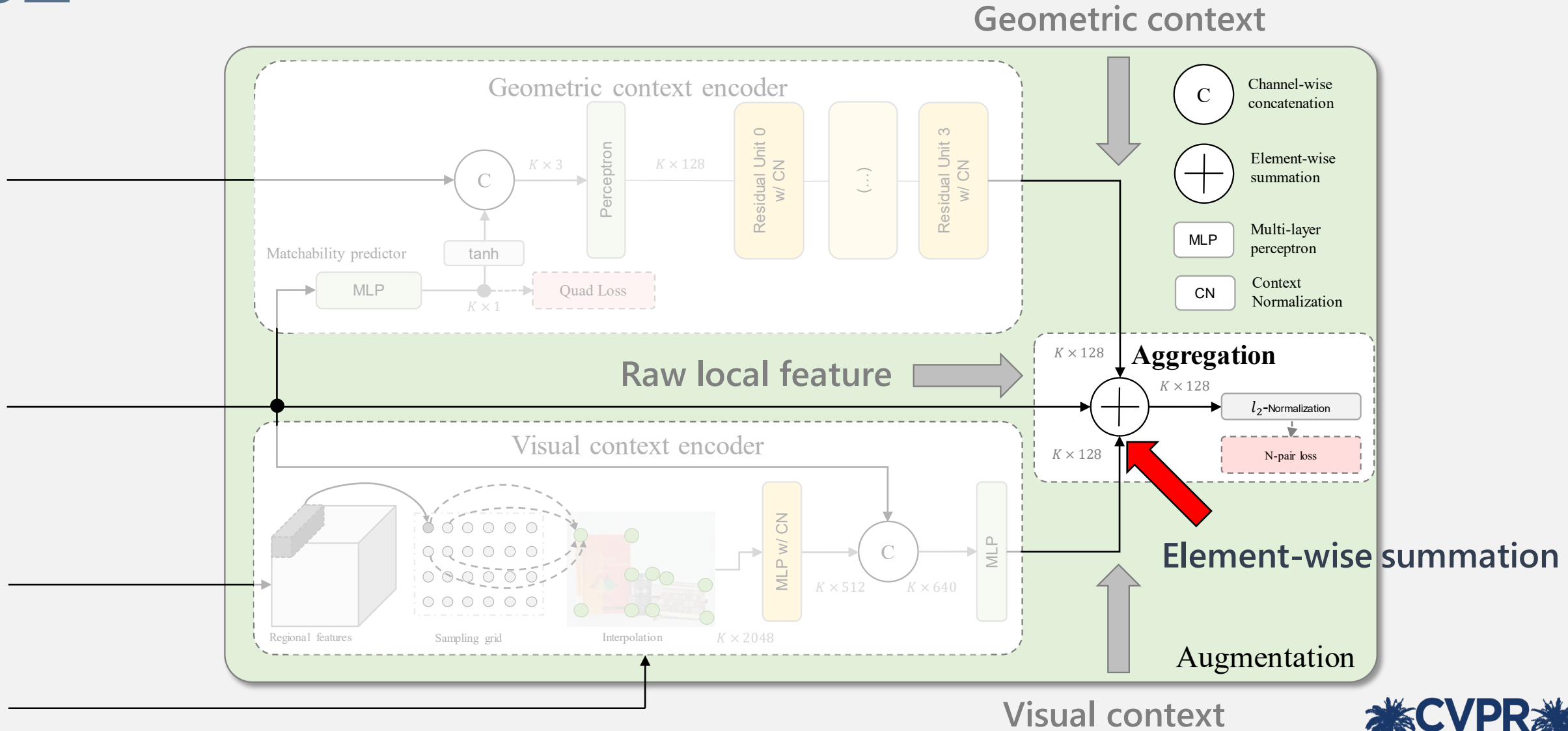
# 02 Methods



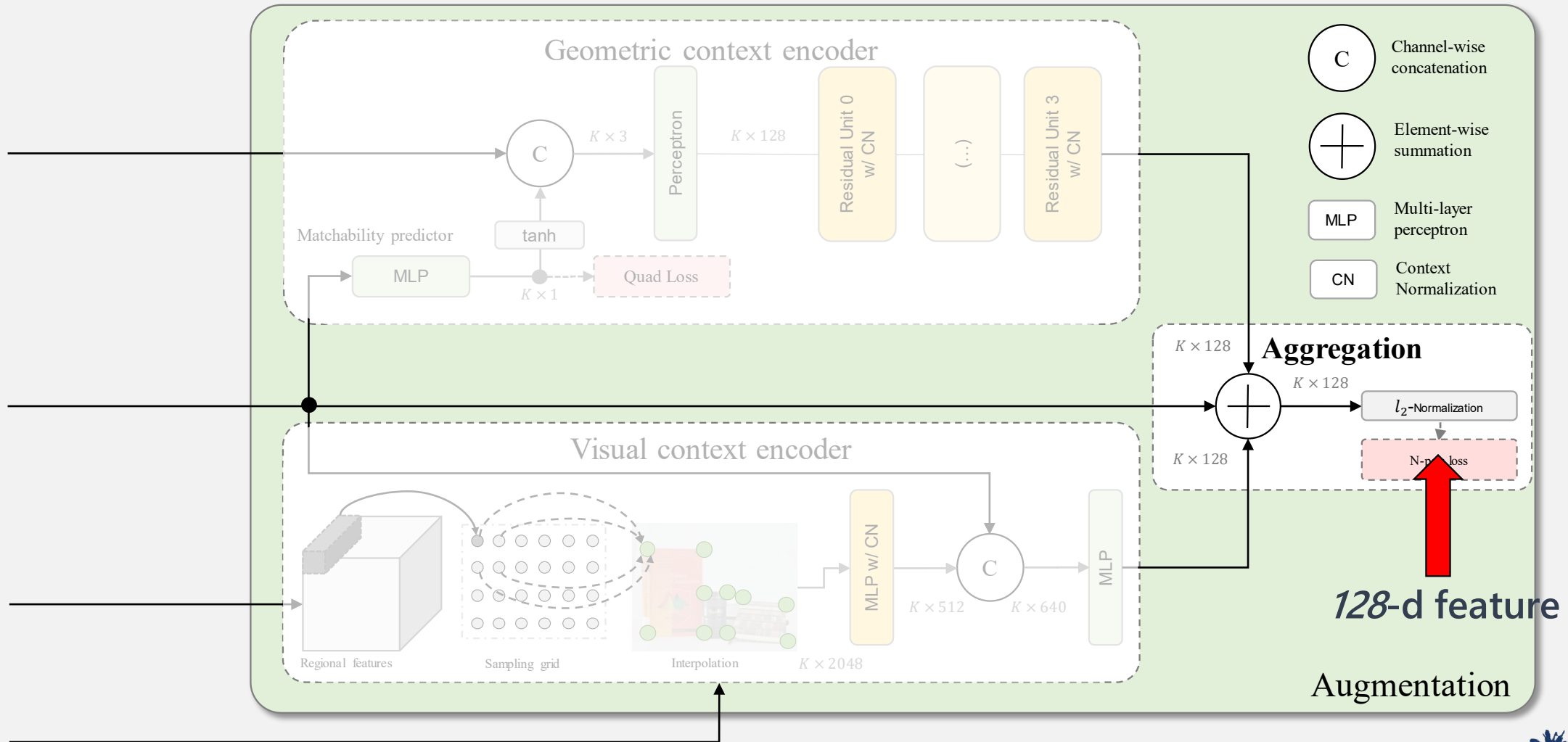
# 02 Methods



# 02 Methods




# 02 Methods



# 03 Ablations

Improvements from  
*visual context*



Visual context encoder		
<i>Strategy</i>	<i>Recall i/v</i>	
baseline (GeoDesc [23])	59.46	71.24
CS (256-d) [50, 19, 43]	59.83	71.27
w/ global feature [5]	59.11	71.02
w/ regional feature	63.64	73.37
<b>w/ regional feature + CN</b>	<b>63.98</b>	<b>73.63</b>

Geometric context encoder		
<i>Network architecture</i>	<i>Recall i/v</i>	
baseline (GeoDesc [23])	59.46	71.24
PointNet [31]	59.61	70.96
w/ CN (pre.) + xy	61.67	72.63
w/ CN (pre.) + xy + raw local feature	60.91	72.99
w/ CN (orig.) + xy + matchability	59.94	71.25
<b>w/ CN (pre.) + xy + matchability</b>	<b>62.82</b>	<b>73.40</b>

Comparison with other methods		
<i>Method</i>	<i>Recall i/v</i>	
SIFT [22]	47.36	53.06
L2-Net [43]	47.58	53.96
HardNet [25]	57.63	63.36
GeoDesc [23]	59.46	71.24
<b>ContextDesc</b>	<b>66.55</b>	<b>75.52</b>
<b>ContextDesc+</b>	<b>67.14</b>	<b>76.42</b>

# 03 Ablations

Improvements from  
*geometric context*



Visual context encoder		
<i>Strategy</i>	<i>Recall i/v</i>	
baseline (GeoDesc [23])	59.46	71.24
CS (256-d) [50, 19, 43]	59.83	71.27
w/ global feature [5]	59.11	71.02
w/ regional feature	63.64	73.37
<b>w/ regional feature + CN</b>	<b>63.98</b>	<b>73.63</b>

Geometric context encoder		
<i>Network architecture</i>	<i>Recall i/v</i>	
baseline (GeoDesc [23])	59.46	71.24
PointNet [31]	59.61	70.96
w/ CN (pre.) + xy	61.67	72.63
w/ CN (pre.) + xy + raw local feature	60.91	72.99
w/ CN (orig.) + xy + matchability	59.94	71.25
<b>w/ CN (pre.) + xy + matchability</b>	<b>62.82</b>	<b>73.40</b>

Comparison with other methods		
<i>Method</i>	<i>Recall i/v</i>	
SIFT [22]	47.36	53.06
L2-Net [43]	47.58	53.96
HardNet [25]	57.63	63.36
GeoDesc [23]	59.46	71.24
<b>ContextDesc</b>	<b>66.55</b>	<b>75.52</b>
<b>ContextDesc+</b>	<b>67.14</b>	<b>76.42</b>

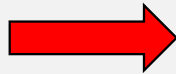


# 03 Ablations

Visual context encoder		
<i>Strategy</i>	<i>Recall i/v</i>	
baseline (GeoDesc [23])	59.46	71.24
CS (256-d) [50, 19, 43]	59.83	71.27
w/ global feature [5]	59.11	71.02
w/ regional feature	63.64	73.37
<b>w/ regional feature + CN</b>	<b>63.98</b>	<b>73.63</b>

Geometric context encoder		
<i>Network architecture</i>	<i>Recall i/v</i>	
baseline (GeoDesc [23])	59.46	71.24
PointNet [31]	59.61	70.96
w/ CN (pre.) + xy	61.67	72.63
w/ CN (pre.) + xy + raw local feature	60.91	72.99
w/ CN (orig.) + xy + matchability	59.94	71.25
<b>w/ CN (pre.) + xy + matchability</b>	<b>62.82</b>	<b>73.40</b>

Improvements from  
*cross-modality context*



Comparison with other methods		
<i>Method</i>	<i>Recall i/v</i>	
SIFT [22]	47.36	53.06
L2-Net [43]	47.58	53.96
HardNet [25]	57.63	63.36
GeoDesc [23]	59.46	71.24
<b>ContextDesc</b>	<b>66.55</b>	<b>75.52</b>
<b>ContextDesc+</b>	<b>67.14</b>	<b>76.42</b>

# 04 Evaluations

SUN3D: indoor scenes

	SIFT [22]	L2-Net [43]	HardNet [25]	GeoDesc [23]	Ours
<i>median number of inlier matches</i>					
<i>indoor</i>	138	153	239	271	<b>365</b>
<i>outdoor</i>	168	173	219	214	<b>482</b>

	SIFT [22]	GeoDesc [23]	Ours
<i>Recall</i>			
<i>JPEG</i>	60.7	66.1	<b>78.6</b>
<i>Blur</i>	41.0	47.7	<b>57.8</b>
<i>Exposure</i>	78.2	86.4	<b>88.2</b>
<i>Day-Night</i>	29.2	39.6	<b>43.3</b>
<i>Scale</i>	81.2	85.8	<b>88.1</b>
<i>Rotation</i>	82.4	<b>87.6</b>	86.3
<i>Scale-Rotation</i>	29.6	33.7	<b>38.0</b>
<i>Planar</i>	48.2	59.1	<b>61.7</b>

		# Images	# Registered	# Sparse Points	# Observations
<b>Fountain</b>	<i>SIFT [22]</i>	11	11	10,004	44K
	<i>GeoDesc [23]</i>		11	16,687	83K
	<i>Ours</i>		11	<b>16,965</b>	<b>84K</b>
<b>Herzjesu</b>	<i>SIFT</i>	8	8	4,916	19K
	<i>GeoDesc</i>		8	8,720	38K
	<i>Ours</i>		8	<b>9,429</b>	<b>40K</b>
<b>South Building</b>	<i>SIFT</i>	128	128	62,780	353K
	<i>GeoDesc</i>		128	170,306	887K
	<i>Ours</i>		128	<b>174,359</b>	<b>893K</b>
<b>Roman Forum</b>	<i>SIFT</i>	2,364	1,407	242,192	1,805K
	<i>GeoDesc</i>		1,566	770,363	5,051K
	<i>Ours</i>		<b>1,571</b>	<b>848,319</b>	<b>5,484K</b>
<b>Alamo</b>	<i>SIFT</i>	2,915	743	120,713	1,384K
	<i>GeoDesc</i>		893	353,329	3,159K
	<i>Ours</i>		<b>921</b>	<b>424,348</b>	<b>3,488K</b>

# 04 Evaluations

YFCC: outdoor scenes

	SIFT [22]	L2-Net [43]	HardNet [5]	GeoDesc [23]	Ours
<i>median number of inlier matches</i>					
<i>indoor</i>	138	153	239	271	<b>365</b>
<i>outdoor</i>	168	173	219	214	<b>482</b>

	SIFT [22]	GeoDesc [23]	Ours
<i>Recall</i>			
<i>JPEG</i>	60.7	66.1	<b>78.6</b>
<i>Blur</i>	41.0	47.7	<b>57.8</b>
<i>Exposure</i>	78.2	86.4	<b>88.2</b>
<i>Day-Night</i>	29.2	39.6	<b>43.3</b>
<i>Scale</i>	81.2	85.8	<b>88.1</b>
<i>Rotation</i>	82.4	<b>87.6</b>	86.3
<i>Scale-Rotation</i>	29.6	33.7	<b>38.0</b>
<i>Planar</i>	48.2	59.1	<b>61.7</b>

		# Images	# Registered	# Sparse Points	# Observations
<b>Fountain</b>	<i>SIFT [22]</i>	11	11	10,004	44K
	<i>GeoDesc [23]</i>		11	16,687	83K
	<i>Ours</i>		11	<b>16,965</b>	<b>84K</b>
<b>Herzjesu</b>	<i>SIFT</i>	8	8	4,916	19K
	<i>GeoDesc</i>		8	8,720	38K
	<i>Ours</i>		8	<b>9,429</b>	<b>40K</b>
<b>South Building</b>	<i>SIFT</i>	128	128	62,780	353K
	<i>GeoDesc</i>		128	170,306	887K
	<i>Ours</i>		128	<b>174,359</b>	<b>893K</b>
<b>Roman Forum</b>	<i>SIFT</i>	2,364	1,407	242,192	1,805K
	<i>GeoDesc</i>		1,566	770,363	5,051K
	<i>Ours</i>		<b>1,571</b>	<b>848,319</b>	<b>5,484K</b>
<b>Alamo</b>	<i>SIFT</i>	2,915	743	120,713	1,384K
	<i>GeoDesc</i>		893	353,329	3,159K
	<i>Ours</i>		<b>921</b>	<b>424,348</b>	<b>3,488K</b>

# 04 Evaluations

	SIFT [22]	L2-Net [43]	HardNet [25]	GeoDesc [23]	Ours
	<i>median number of inlier matches</i>				
<i>indoor</i>	138	153	239	271	<b>365</b>
<i>outdoor</i>	168	173	219	214	<b>482</b>

	SIFT [22]	GeoDesc [23]	Ours
	<i>Recall</i>		
<i>JPEG</i>	60.7	66.1	<b>78.6</b>
<i>Blur</i>	41.0	47.7	<b>57.8</b>
<i>Exposure</i>	78.2	86.4	<b>88.2</b>
<i>Day-Night</i>	29.2	39.6	<b>43.3</b>
<i>Scale</i>	81.2	85.8	<b>88.1</b>
<i>Rotation</i>	82.4	<b>87.6</b>	86.3
<i>Scale-Rotation</i>	29.6	33.7	<b>38.0</b>
<i>Planar</i>	48.2	59.1	<b>61.7</b>



Oxford dataset:  
Different image  
variations

		# Images	# Registered	# Sparse Points	# Observations
<b>Fountain</b>	<i>SIFT [22]</i>	11	11	10,004	44K
	<i>GeoDesc [23]</i>		11	16,687	83K
	<i>Ours</i>		11	<b>16,965</b>	<b>84K</b>
<b>Herzjesu</b>	<i>SIFT</i>	8	8	4,916	19K
	<i>GeoDesc</i>		8	8,720	38K
	<i>Ours</i>		8	<b>9,429</b>	<b>40K</b>
<b>South Building</b>	<i>SIFT</i>	128	128	62,780	353K
	<i>GeoDesc</i>		128	170,306	887K
	<i>Ours</i>		128	<b>174,359</b>	<b>893K</b>
<b>Roman Forum</b>	<i>SIFT</i>	2,364	1,407	242,192	1,805K
	<i>GeoDesc</i>		1,566	770,363	5,051K
	<i>Ours</i>		<b>1,571</b>	<b>848,319</b>	<b>5,484K</b>
<b>Alamo</b>	<i>SIFT</i>	2,915	743	120,713	1,384K
	<i>GeoDesc</i>		893	353,329	3,159K
	<i>Ours</i>		<b>921</b>	<b>424,348</b>	<b>3,488K</b>

# 04 Evaluations

	SIFT [22]	L2-Net [43]	HardNet [25]	GeoDesc [23]	Ours
<i>median number of inlier matches</i>					
<i>indoor</i>	138	153	239	271	<b>365</b>
<i>outdoor</i>	168	173	219	214	<b>482</b>

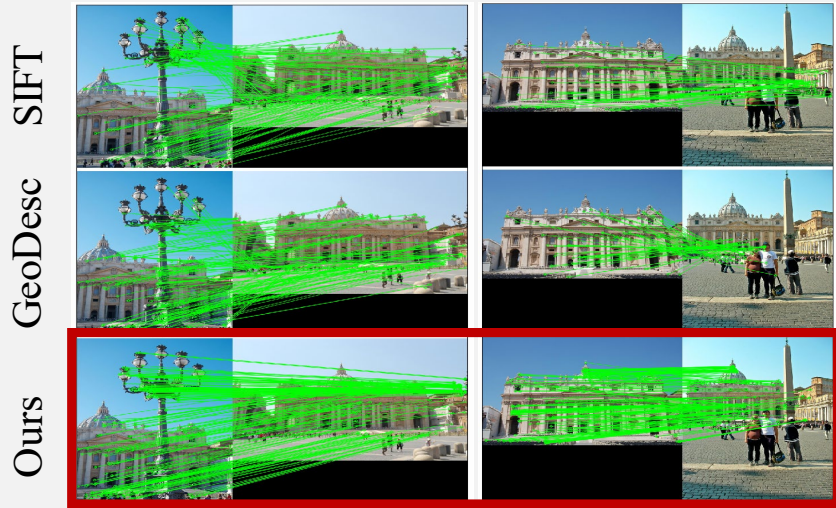
	SIFT [22]	GeoDesc [23]	Ours
<i>Recall</i>			
<i>JPEG</i>	60.7	66.1	<b>78.6</b>
<i>Blur</i>	41.0	47.7	<b>57.8</b>
<i>Exposure</i>	78.2	86.4	<b>88.2</b>
<i>Day-Night</i>	29.2	39.6	<b>43.3</b>
<i>Scale</i>	81.2	85.8	<b>88.1</b>
<i>Rotation</i>	82.4	<b>87.6</b>	86.3
<i>Scale-Rotation</i>	29.6	33.7	<b>38.0</b>
<i>Planar</i>	48.2	59.1	<b>61.7</b>

		# Images	# Registered	# Sparse Points	# Observations
<b>Fountain</b>	<i>SIFT [22]</i>	11	11	10,004	44K
	<i>GeoDesc [23]</i>		11	16,687	83K
	<i>Ours</i>		11	<b>16,965</b>	<b>84K</b>
<b>Herzjesu</b>	<i>SIFT</i>	8	8	4,916	19K
	<i>GeoDesc</i>		8	8,720	38K
	<i>Ours</i>		8	<b>9,429</b>	<b>40K</b>
<b>South Building</b>	<i>SIFT</i>	128	128	62,780	353K
	<i>GeoDesc</i>		128	170,306	887K
	<i>Ours</i>		128	<b>174,359</b>	<b>893K</b>
<b>Roman Forum</b>	<i>SIFT</i>	2,364	1,407	242,192	1,805K
	<i>GeoDesc</i>		1,566	770,363	5,051K
	<i>Ours</i>		<b>1,571</b>	<b>848,319</b>	<b>5,484K</b>
<b>Alamo</b>	<i>SIFT</i>	2,915	743	120,713	1,384K
	<i>GeoDesc</i>		893	353,329	3,159K
	<i>Ours</i>		<b>921</b>	<b>424,348</b>	<b>3,488K</b>

3D reconstruction  
benchmark



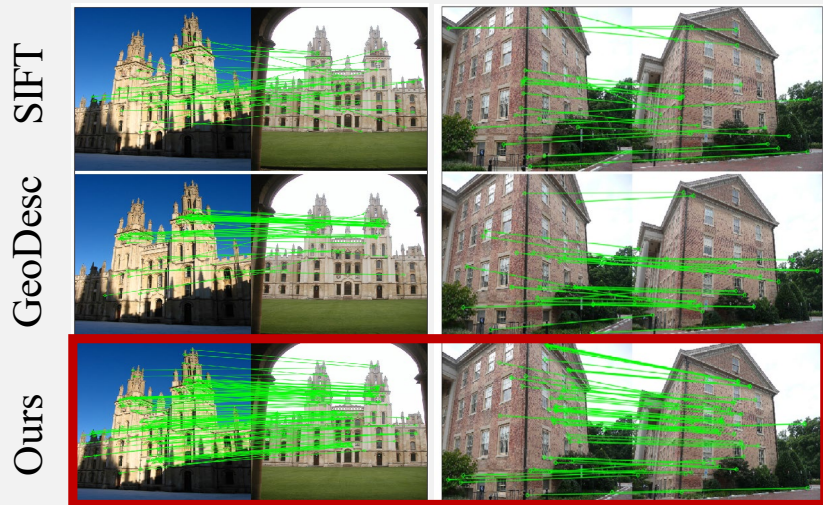
# 04 Evaluations



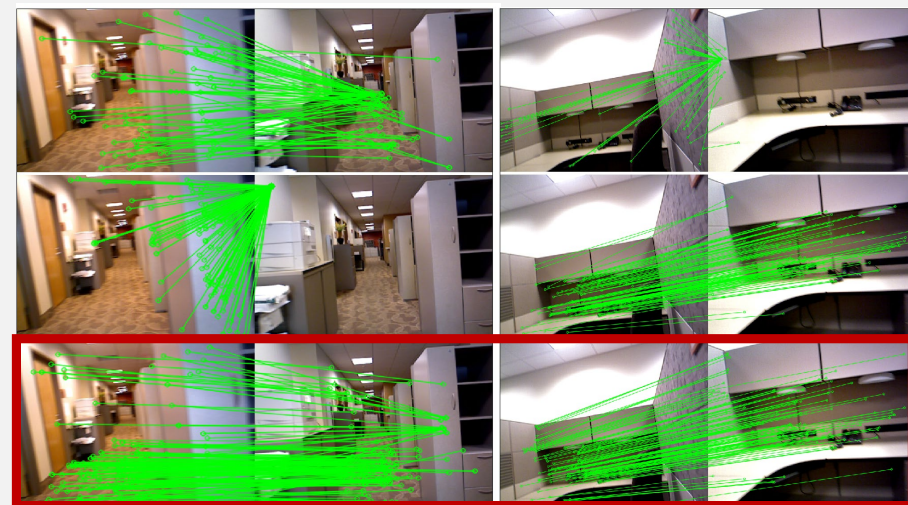
Scale or rotation change



Illumination change





Perspective change



Indoor scene

# 04 Evaluations

When pose accuracy as evaluation metric: consistent improvement, but less significant

Method	Date	Type	Ims (%)	#Pts	SR	TL	mAP <sup>5°</sup>	mAP <sup>10°</sup>	mAP <sup>15°</sup>	mAP <sup>20°</sup>	mAP <sup>25°</sup>	ATE
 SIFT + ContextDesc kp:8000, match:nn	19-05-09	F	97.9	6020.4	97.2	3.26	0.3828	0.4821	0.5399	0.5853	0.6226	—
 SIFT + GeoDesc kp:8000, match:nn	19-04-24	F	97.3	5583.8	95.8	3.39	0.3858	0.4778	0.5317	0.5790	0.6139	—

After obtaining sufficient matches, what is the next bottleneck in order to improve the image matching?

# 05 Sparse matching

- 1) Establish putative matches (nearest-neighbor search/FLANN)
- 2) Outlier rejection (ratio test/mutual check/GMS)
- 3) Geometry computation (5-point/8-point algorithm with RANSAC)
- 4) Non-linear optimization for refinement



# 05 Sparse matching

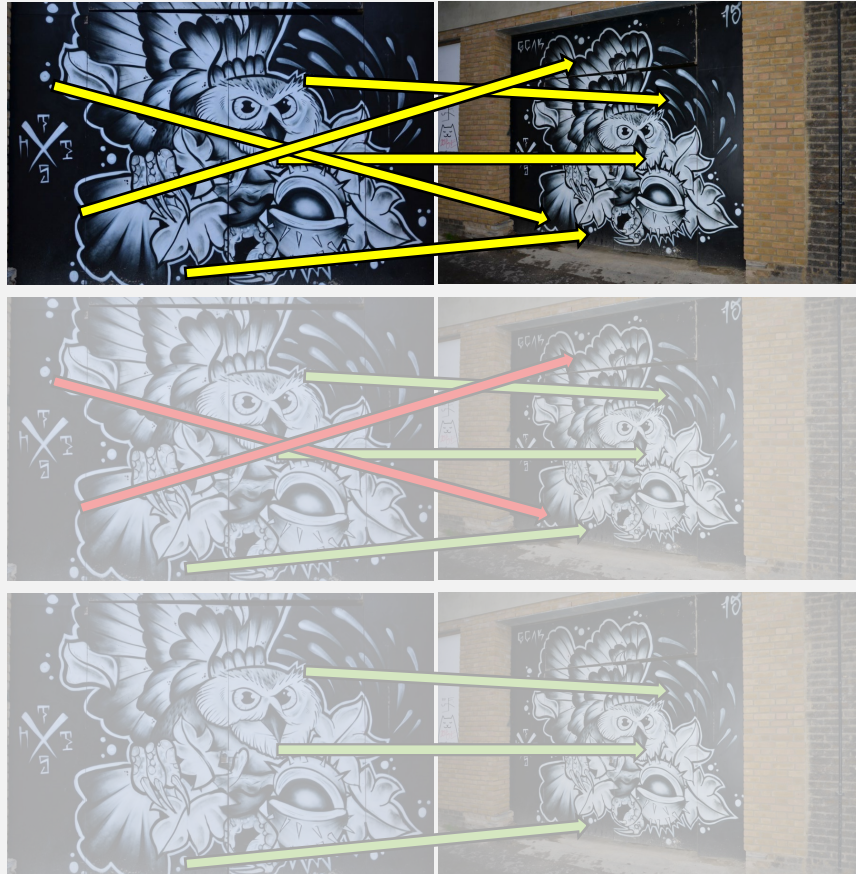
- 1) Establish putative matches (nearest-neighbor search/FLANN)
- 2) **Outlier rejection (ratio test/mutual check/GMS)**
- 3) Geometry computation (5-point/8-point algorithm with RANSAC)
- 4) Non-linear optimization for refinement

Learning-based

\*Yi et al.: Learning to find good correspondences, CVPR'18.



# 05 Sparse matching

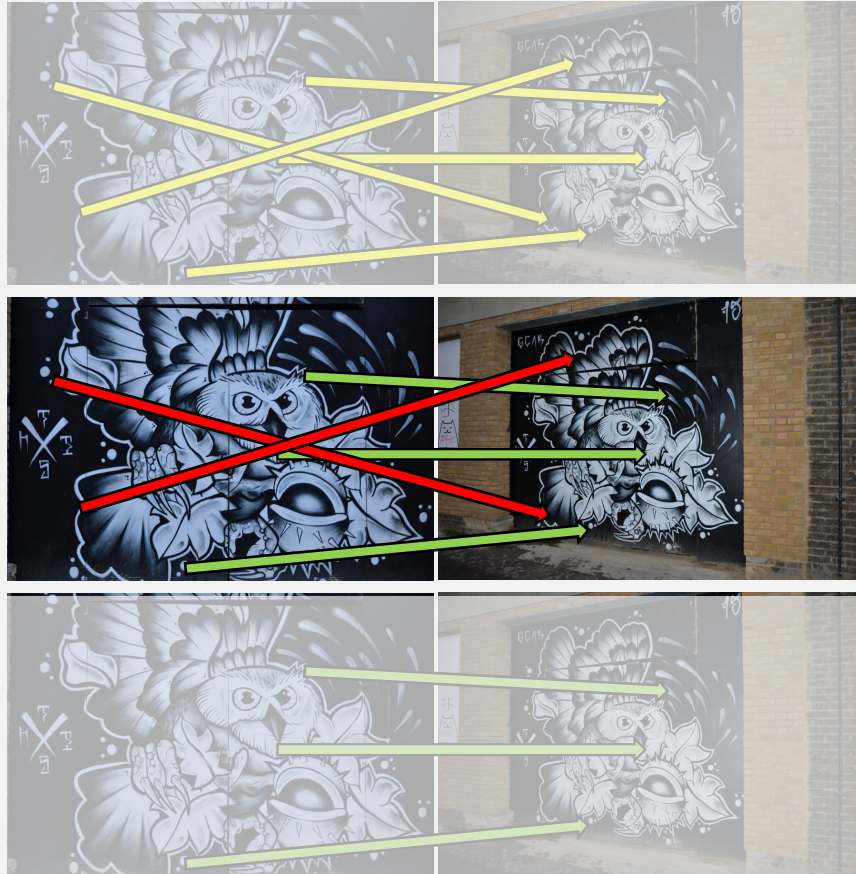


Given putative matches  $N \times 4$ , where each row vector denotes a correspondence  $(x, y, x', y')$  of an image pair.

The network predicts the probability vector  $N \times 1$  that indicates whether a correspondence is an inlier.

Only inlier matches (and its confidence) are used for computing the geometry.

# 05 Sparse matching



Given putative matches  $N \times 4$ , where each row vector denotes a correspondence  $(x, y, x', y')$  of an image pair.

The network predicts the probability vector  $N \times 1$  that indicates whether a correspondence is an inlier.

Only inlier matches (and its confidence) are used for computing the geometry.

# 05 Sparse matching



Given putative matches  $N \times 4$ , where each row vector denotes a correspondence  $(x, y, x', y')$  of an image pair.

The network predicts the probability vector  $N \times 1$  that indicates whether a correspondence is an inlier.

Only inlier matches (and their confidence) are used for solving the two-view geometry.

# 05 Sparse matching

Why is it important?

Method	Date	Type	Ims (%)	#Pts	SR	TL	mAP <sup>5°</sup>	mAP <sup>10°</sup>	mAP <sup>15°</sup>	mAP <sup>20°</sup>	mAP <sup>25°</sup>	ATE
<b>SIFT + ContextDesc + Inlier Classification V2</b> kp:8000, match:custom	19-05-28	F	98.1	6126.0	97.5	3.44	0.5755	0.6830	0.7389	0.7750	0.8006	—
<b>SIFT + ContextDesc</b> kp:8000, match:nn1to1	19-06-07	F	98.1	6472.1	98.0	3.34	0.4287	0.5371	0.6017	0.6464	0.6826	—
<b>SIFT + ContextDesc</b> kp:8000, match:nn	19-05-09	F	97.9	6020.4	97.2	3.26	0.3828	0.4821	0.5399	0.5853	0.6226	—

Proposed learning-based

Mutual check

No outlier rejection

# 05 Sparse matching

## Previous method

- Adopt a PointNet-like architecture.
- Apply context normalization (instance normalization) on the entire point set to capture **global** context.

Local context, e.g., piece-wise smoothness (GMS matcher).



\*Yi et al.: Learning to find good correspondences, CVPR'18.

# 05 Sparse matching

## Previous method

- Adopt a PointNet-like architecture.
- Apply context normalization (instance normalization) on the entire point set to capture **global context**.

Local context, e.g., piece-wise smoothness (GMS matcher).



\*Bian et al.: GMS: Grid-based Motion Statistics for Fast, Ultra-robust Feature Correspondence, CVPR'17.

# 05 Sparse matching

## Previous method

- Adopt a PointNet-like architecture.
- Apply context normalization (instance normalization) on the entire point set to capture **global** context.

## Proposed

- Learn to establish neighboring relations on unordered, non-Euclidean correspondence sets.
- Build a **hierarchical** architecture to capture both **global and local context**.



# 05 Sparse matching

Method	Date	Type	Ims (%)	#Pts	SR	TL	mAP <sup>5°</sup>	mAP <sup>10°</sup>	mAP <sup>15°</sup>	mAP <sup>20°</sup>	mAP <sup>25°</sup>	ATE
<b>SIFT + ContextDesc + Inlier Classification V2</b> kp:8000, match:custom <b>Proposed</b>	19-05-28	F/M	98.6	6126.0	97.5	3.44	0.5755	0.6830	0.7389	0.7750	0.8006	—
<b>SIFT + ContextDesc + Inlier Classification V1</b> kp:8000, match:custom <b>Previous method</b>	19-05-29	F/M	98.4	6045.8	97.8	3.43	0.5553	0.6633	0.7169	0.7545	0.7849	—

# 06 Future work

## Evaluation metric

- HPatches: patch verification/matching/retrieval - *Reflect the performance in real applications? [1]*
- Two-view image matching on pose recovery accuracy - *Involve other variables such as RANSAC-based algorithms? [2]*
- 3D reconstruction metrics – *Involve more variables such as image retrieval, SfM or bundle adjustment? [3]*
- The ground truth is often obtained from SfM with a traditional matching pipeline.

[1] Balntas et al.: HPatches: A benchmark and evaluation of handcrafted and learned local descriptors, CVPR'17

[2] Yi et al.: Learning to find good correspondences, CVPR'18.

[3] Schönberger et al.: Comparative Evaluation of Hand-Crafted and Learned Local Features, CVPR'17.

## Keypoint detection

- Challenging for learning-based methods – clear definition as supervision?
- Pixel-wise, even sub-pixel wise accuracy – preserve low-level details after multiple convolutions?
- Combine the advantage of both human priors and learned priors.



# 06 Future work

## Evaluation metric

- HPatches: patch verification/matching/retrieval - *Reflect the performance in real applications? [1]*
- Two-view image matching on pose recovery accuracy - *Involve other variables such as RANSAC-based algorithms? [2]*
- 3D reconstruction metrics – *Involve more variables such as image retrieval, SfM or bundle adjustment? [3]*
- The ground truth is often obtained from SfM with a traditional matching pipeline.

[1] Balntas et al.: HPatches: A benchmark and evaluation of handcrafted and learned local descriptors, CVPR'17

[2] Yi et al.: Learning to find good correspondences, CVPR'18.

[3] Schönberger et al.: Comparative Evaluation of Hand-Crafted and Learned Local Features, CVPR'17.

## Keypoint detection

- Challenging for learning-based methods – clear definition as supervision?
- Pixel-wise, even sub-pixel wise accuracy – preserve low-level details after multiple convolutions?
- Combine the advantage of both human priors and learned priors.



# 06 Future work

## Evaluation metric

- HPatches: patch verification/matching/retrieval - *Reflect the performance in real applications? [1]*
- Two-view image matching on pose recovery accuracy - *Involve other variables such as RANSAC-based algorithms? [2]*
- 3D reconstruction metrics – *Involve more variables such as image retrieval, SfM or bundle adjustment? [3]*
- **The ground truth is often obtained from SfM with a traditional matching pipeline.**

[1] Balntas et al.: HPatches: A benchmark and evaluation of handcrafted and learned local descriptors, CVPR'17

[2] Yi et al.: Learning to find good correspondences, CVPR'18.

[3] Schönberger et al.: Comparative Evaluation of Hand-Crafted and Learned Local Features, CVPR'17.

## Keypoint detection

- Challenging for learning-based methods – clear definition as supervision?
- Pixel-wise, even sub-pixel wise accuracy – preserve low-level details after multiple convolutions?
- Combine the advantage of both human priors and learned priors.



# 06 Future work

## Evaluation metric

- HPatches: patch verification/matching/retrieval - *Reflect the performance in real applications? [1]*
- Two-view image matching on pose recovery accuracy - *Involve other variables such as RANSAC-based algorithms? [2]*
- 3D reconstruction metrics – *Involve more variables such as image retrieval, SfM or bundle adjustment? [3]*
- The ground truth is often obtained from SfM with a traditional matching pipeline.

[1] Balntas et al.: HPatches: A benchmark and evaluation of handcrafted and learned local descriptors, CVPR'17

[2] Yi et al.: Learning to find good correspondences, CVPR'18.

[3] Schönberger et al.: Comparative Evaluation of Hand-Crafted and Learned Local Features, CVPR'17.

## Keypoint detection

- Challenging for learning-based methods – clear definition as supervision?
- Pixel-wise, even sub-pixel wise accuracy – preserve low-level details after multiple convolutions?
- Combine the advantage of both human priors and learned priors.



# 06 Future work

## Evaluation metric

- HPatches: patch verification/matching/retrieval - *Reflect the performance in real applications? [1]*
- Two-view image matching on pose recovery accuracy - *Involve other variables such as RANSAC-based algorithms? [2]*
- 3D reconstruction metrics – *Involve more variables such as image retrieval, SfM or bundle adjustment? [3]*
- The ground truth is often obtained from SfM with a traditional matching pipeline.

[1] Balntas et al.: HPatches: A benchmark and evaluation of handcrafted and learned local descriptors, CVPR'17

[2] Yi et al.: Learning to find good correspondences, CVPR'18.

[3] Schönberger et al.: Comparative Evaluation of Hand-Crafted and Learned Local Features, CVPR'17.

## Keypoint detection

- Challenging for learning-based methods – clear definition as supervision?
- Pixel-wise, even sub-pixel wise accuracy – preserve low-level details after multiple convolutions?
- **Combine the advantage of both human priors and learned priors.**



---

# Thanks!

---

Code available at: <https://github.com/lzx551402/contextdesc>



**LONG BEACH**  
**CALIFORNIA**  
**June 16-20, 2019**