

Outdoor Visual Place Recognition Based on 3D Point Cloud

Supervisor: Prof. TEO Chee Leong

Examiner: Prof. TAY Eng Hock

YUAN Chengran

Contents

I. Introduction

II. Clarification/ Problem Definition

III. Literature review

IV. Methodology

V. Experiments

VI. Conclusion

VII. Future work

I. Introduction

Visual Place Recognition (VPR)

What is VPR?

- For an agent (a robot or a vehicle), the ability to recognize the same place despite significant changes in appearance and viewpoints.



Two pictures taken at the same place
(different seasons, weather and illumination conditions)

2D & 3D Methods for VPR

- Based on the data input, there are mainly two approaches to solving VPR,
 - 2D method using **images** as input
 - 3D method using **3D data** (usually **point cloud**) as input.

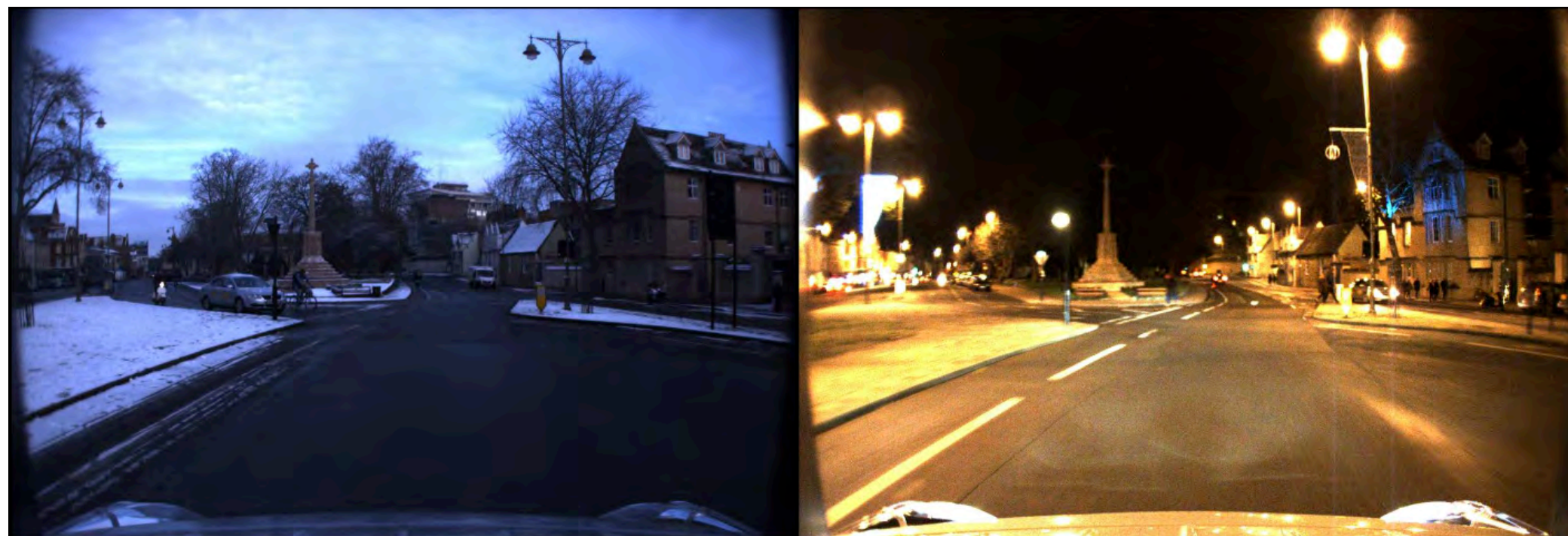
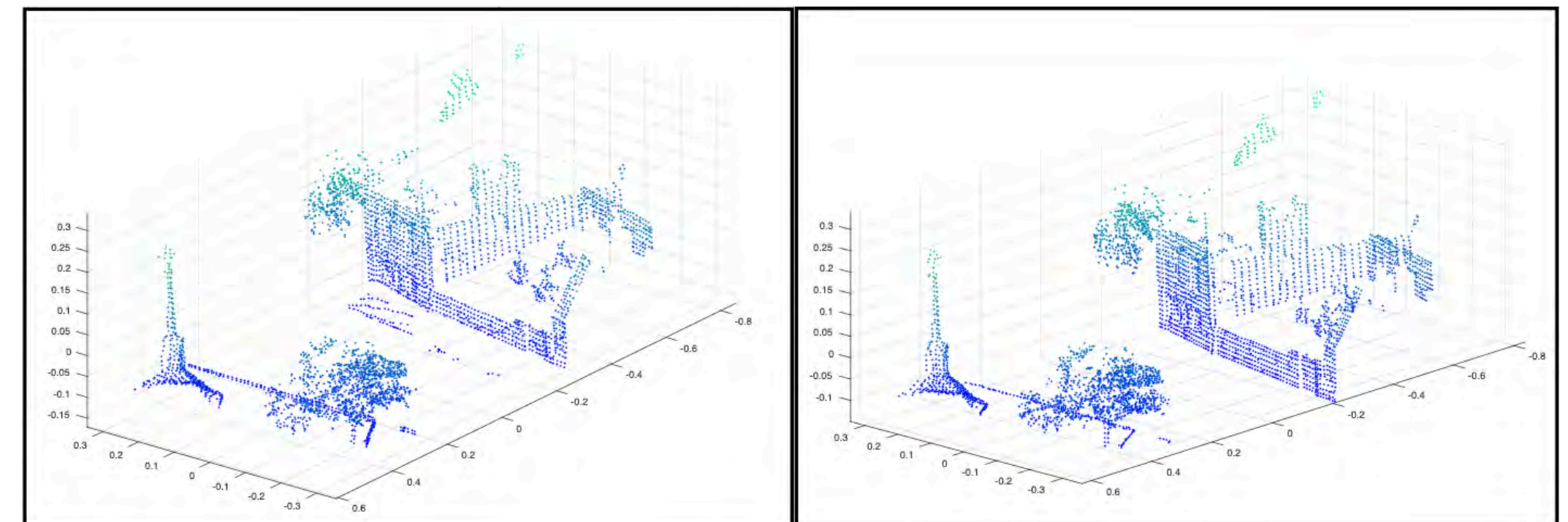


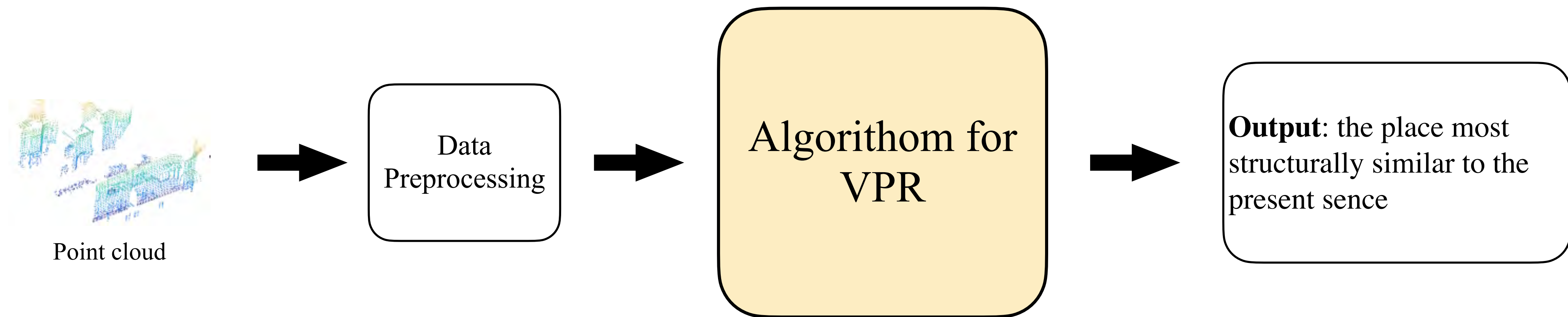
Image input



Point Cloud input

II. Clarification

General framework for 3D VPR



Clarification

The problem defined as follows:

Given a query 3D point cloud denoted as q , where

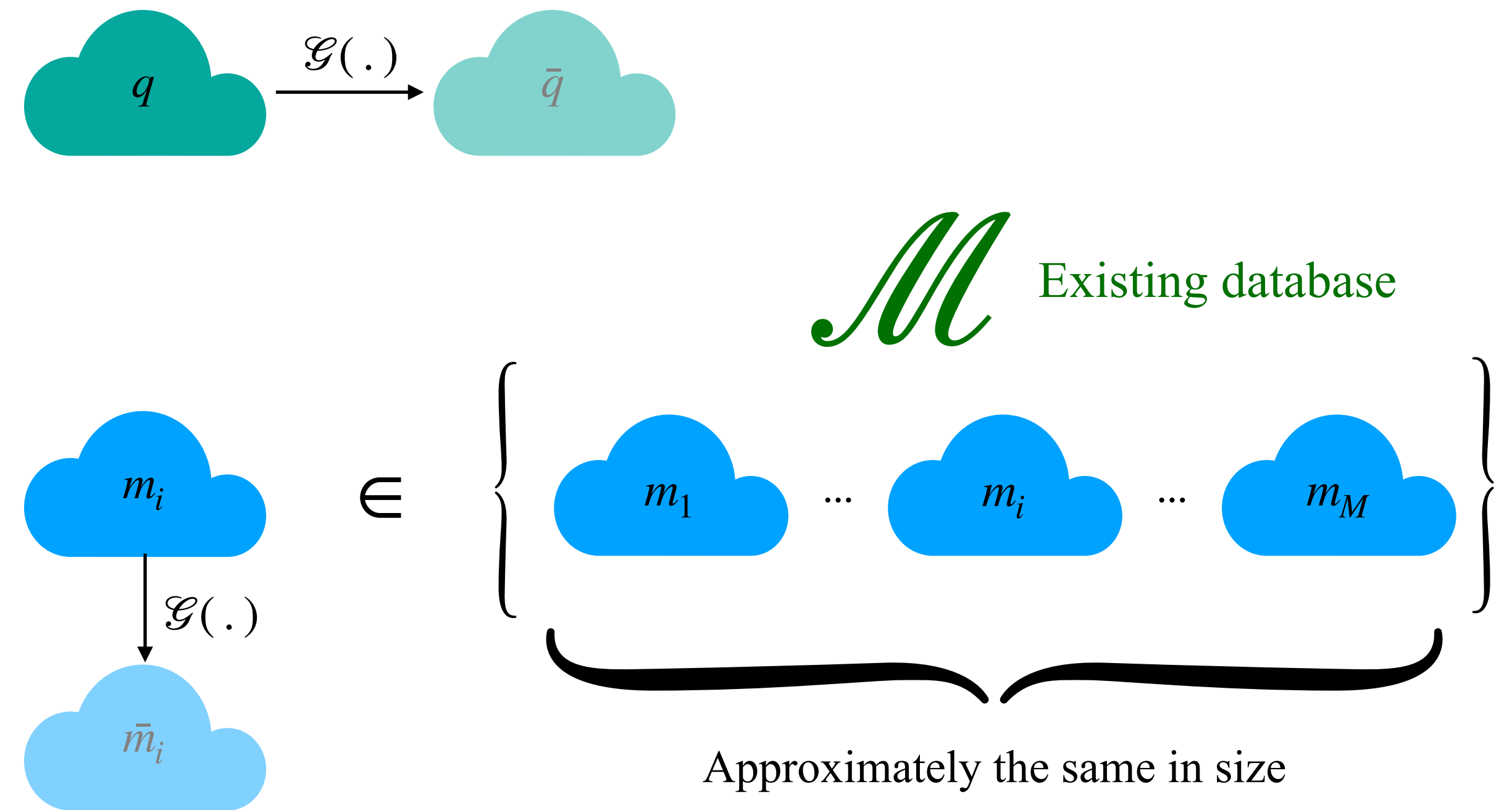
$$AOC(q) \approx AOC(m_i)$$

and

$$|\mathcal{G}(q)| = |\mathcal{G}(m_i)|,$$

AOC: area of coverage $\mathcal{G}(.)$: downsampling filter

The goal is to retrieve the submap m_* from the database \mathcal{M} that is structurally most similar to q .



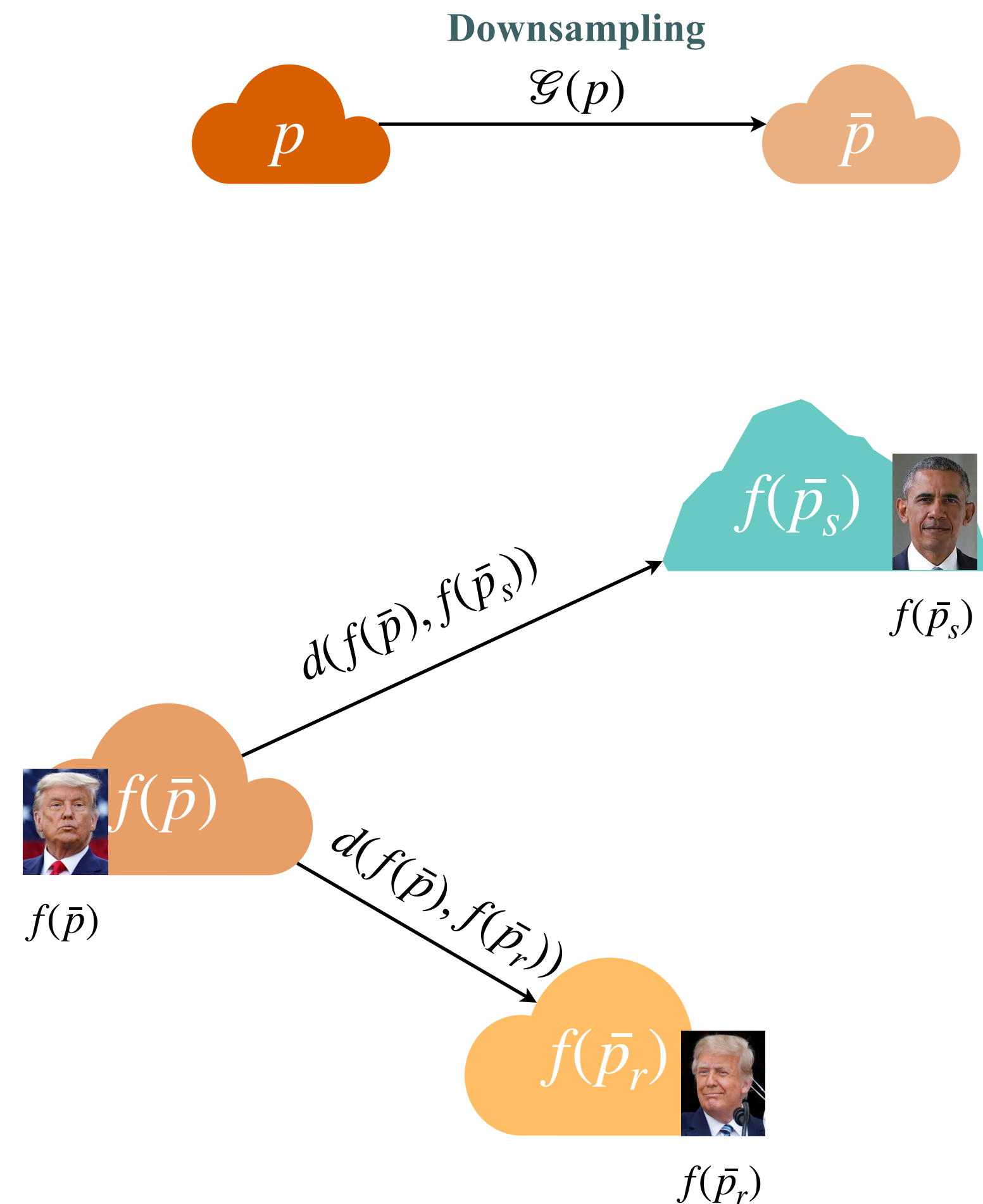
Function Definition

Towards this goal, a deep network is devised to learn a function $f(.)$ that maps a given downsampled 3D point cloud $\bar{p} = \mathcal{G}(p)$ to a fixed size **global descriptor vector** $f(\bar{p})$ such that

$$d(f(\bar{p}), f(\bar{p}_r)) < d(f(\bar{p}), f(\bar{p}_s)),$$

if p is structurally similar to p_r but dissimilar to p_s .

$d(.)$ is some distance function, e.g. Euclidean distance function.



Simplification

Our problem then simplifies to finding the submap

$$m_* \in \mathcal{M}$$

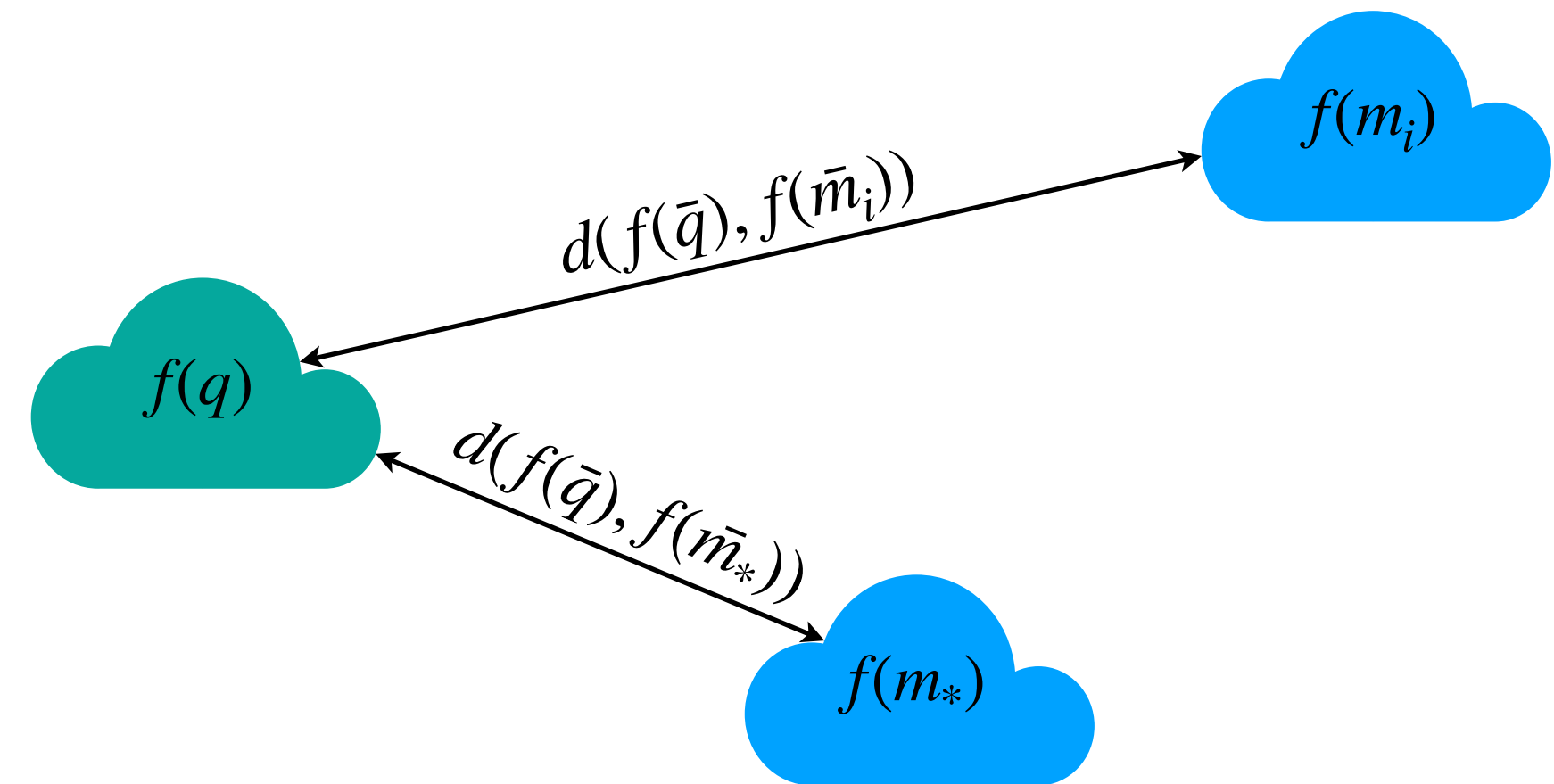
such that its global descriptor vector $f(\bar{m}_*)$ gives the **minimum distance** with the global descriptor vector $f(\bar{q})$ from the query q , i.e.

$$d(f(\bar{q}), f(\bar{m}_*)) < d(f(\bar{q}), f(\bar{m}_i)), \forall i \neq *.$$

In practice, this can be done by **the nearest neighbor search** through a list of global descriptors

$$\{ f(\bar{m}_i) \mid i \in 1, 2, \dots, M \}$$

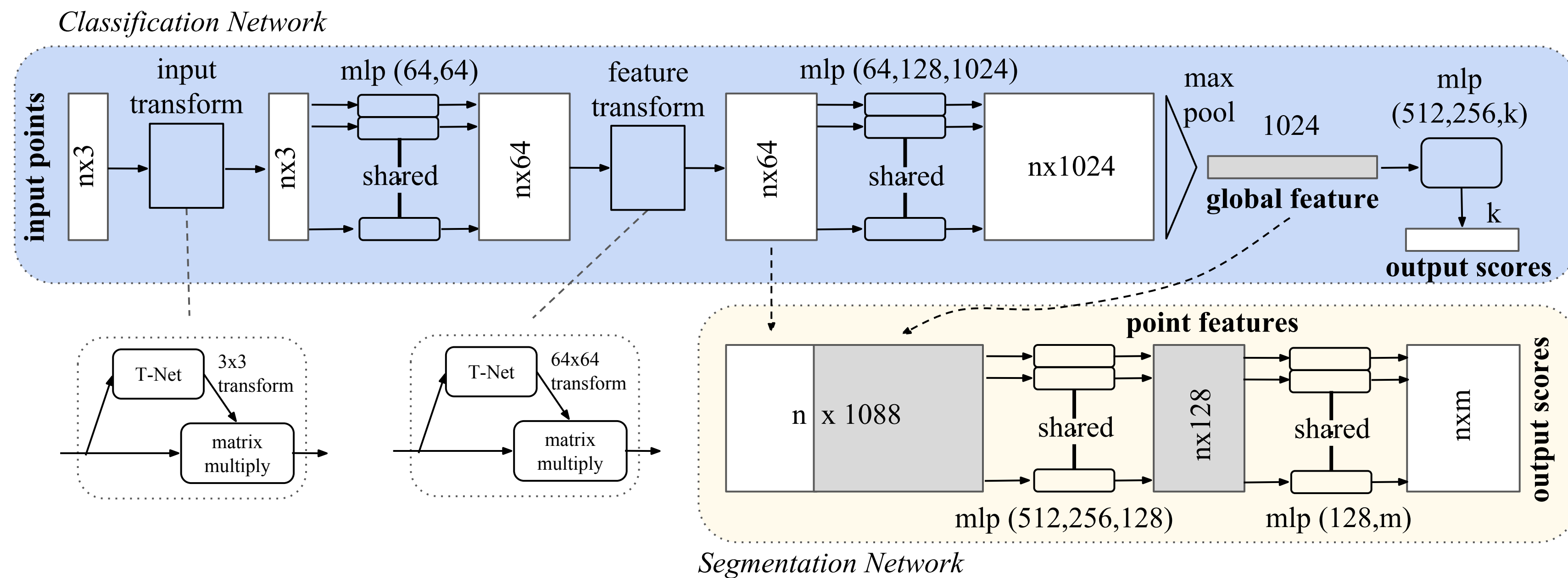
that can be computed once **offline** and stored in memory, while $f(\bar{q})$ is computed **online**.



III. Literature review

PointNet

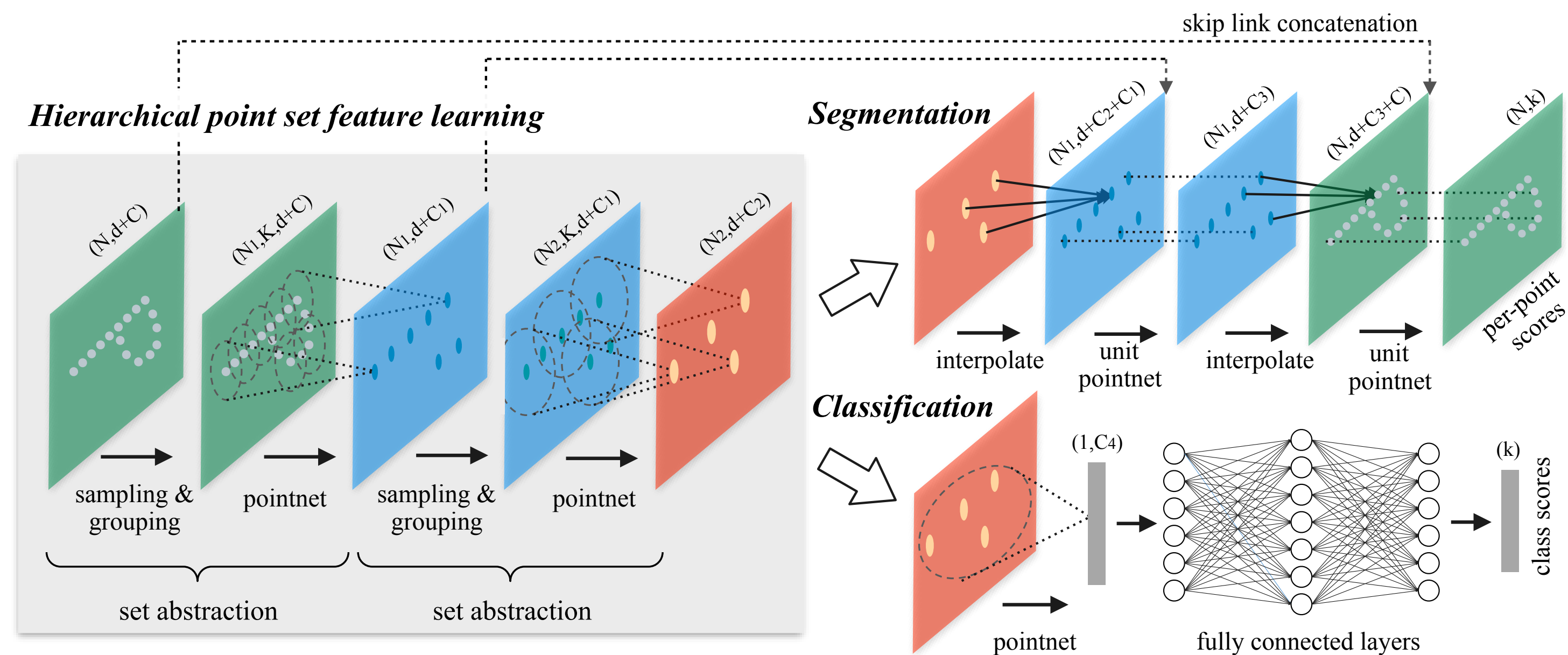
Pioneer in Point Cloud Processing



Qi, C. et al. "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation." *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017): 77-85.

PointNet++

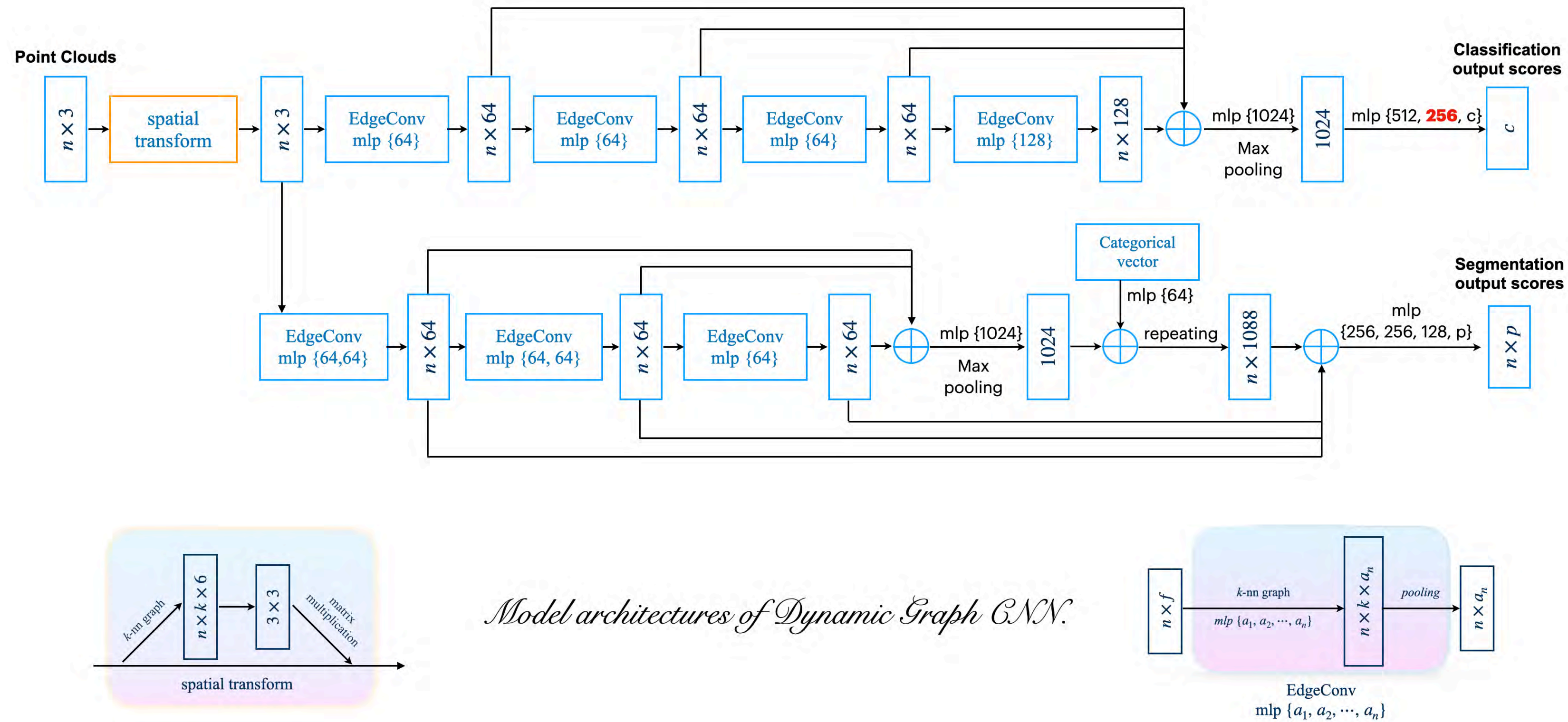
Hierarchical architecture

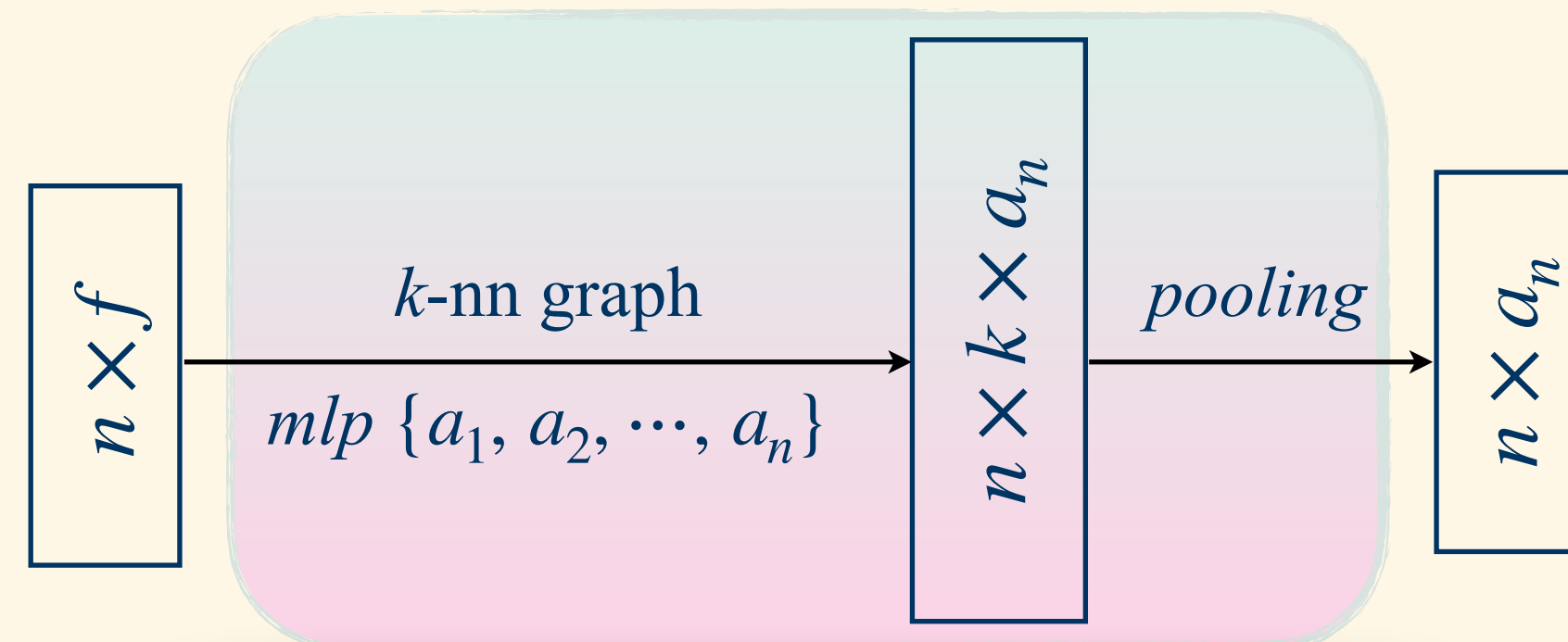


PointNet++: Hierarchical feature learning architecture and its application.

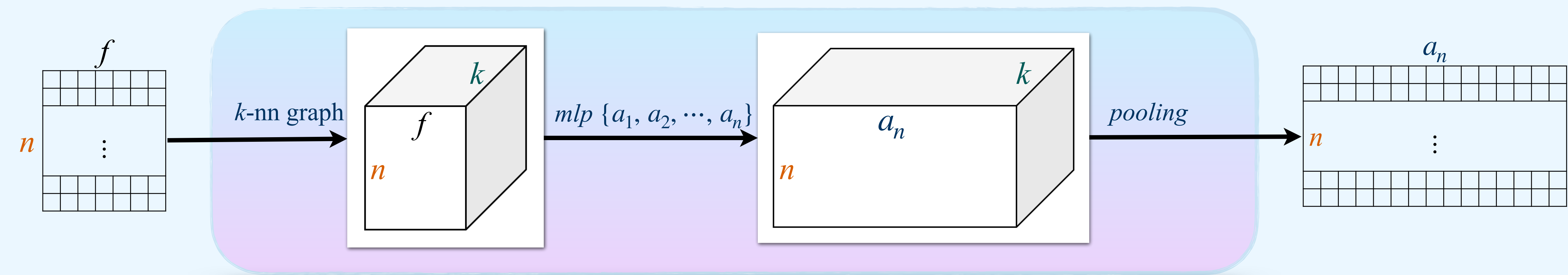
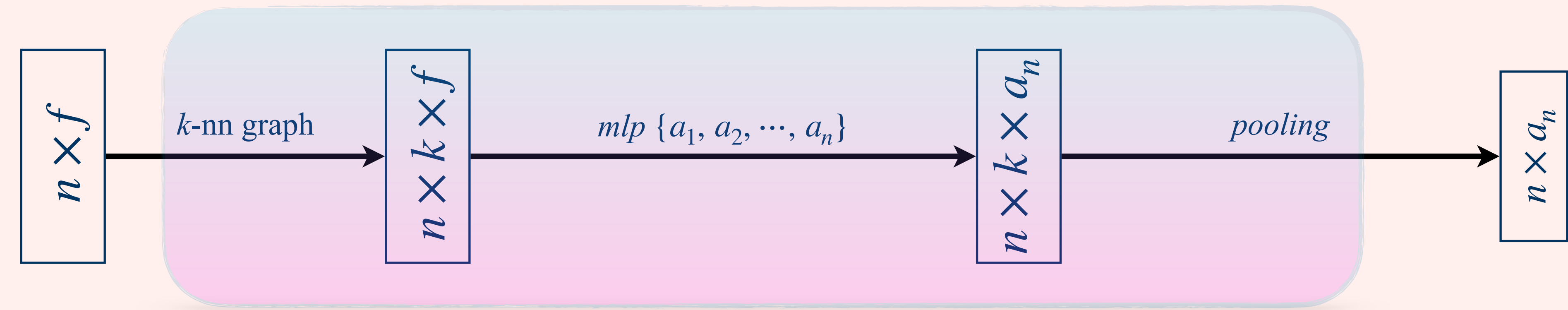
DGCNN

Better extract local geometric features



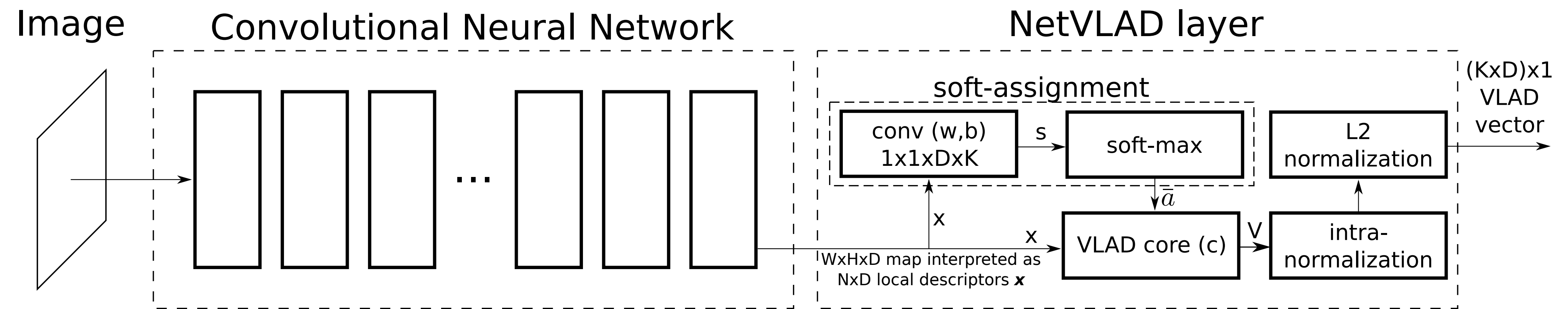


EdgeConv
 $mlp \{a_1, a_2, \dots, a_n\}$



NetVLAD

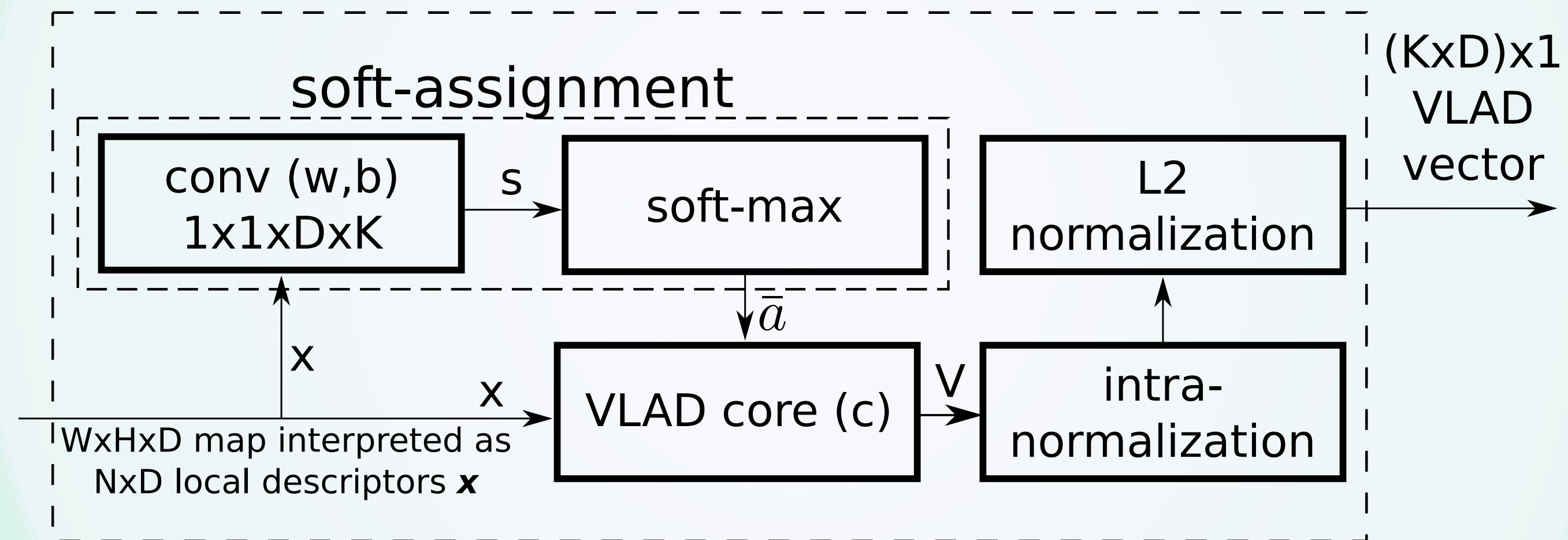
A Backbone for 2D VPR



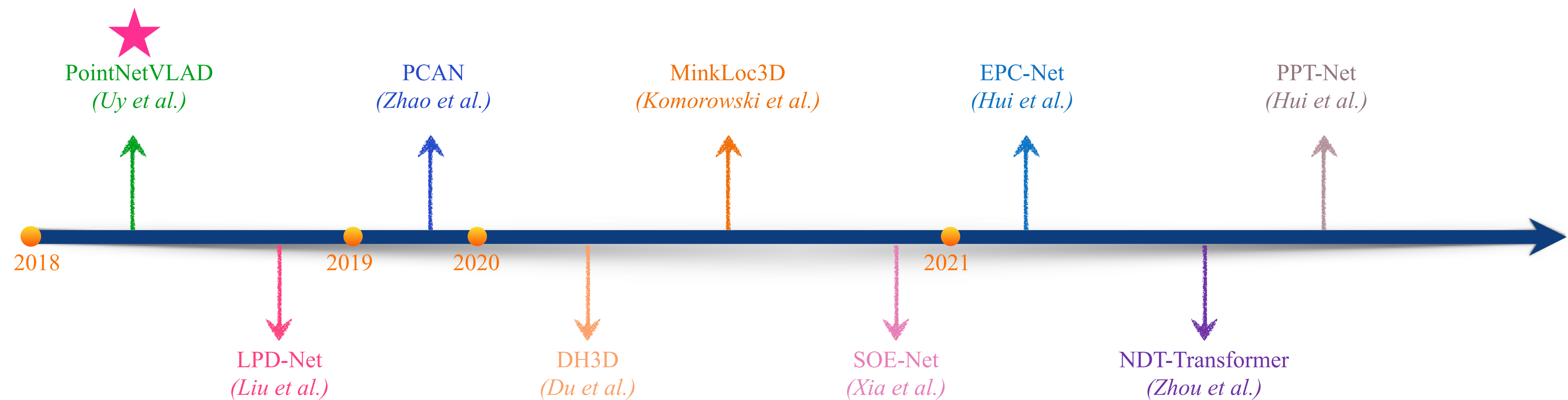
CNN architecture with the NetVLAD layer.

Aggregate local descriptors

NetVLAD layer

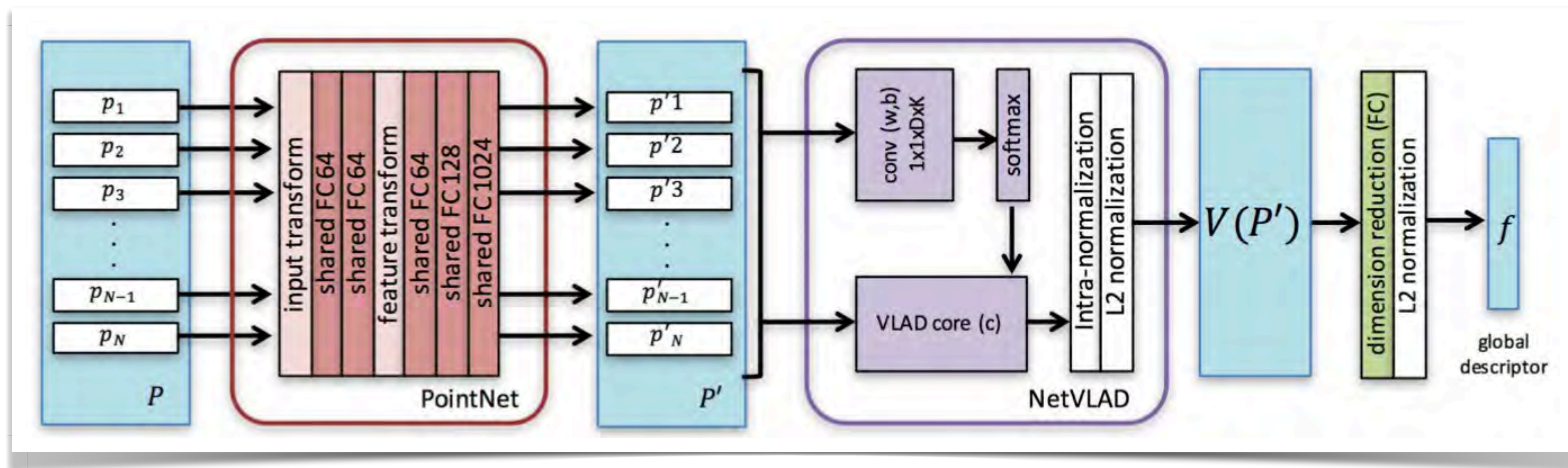


Chronological overview of 3D VPR



PointNetVLAD

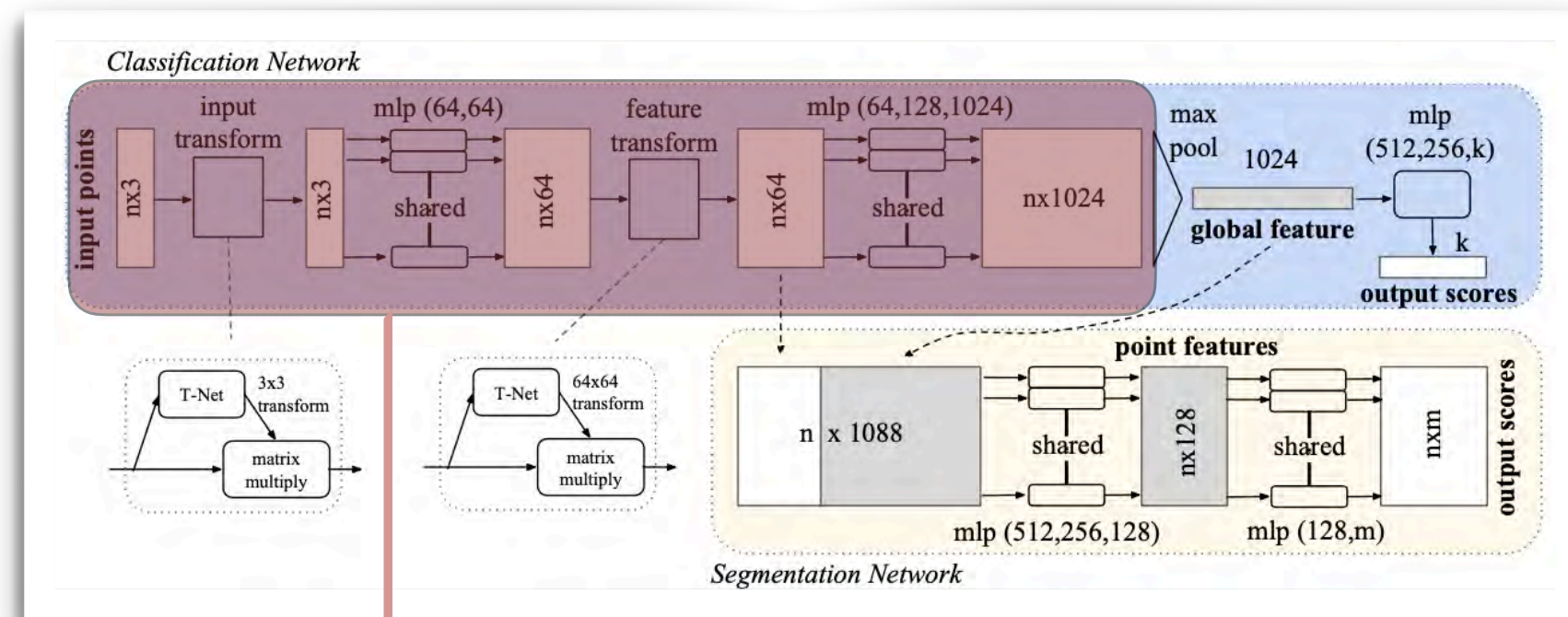
Pioneer of 3D VPR



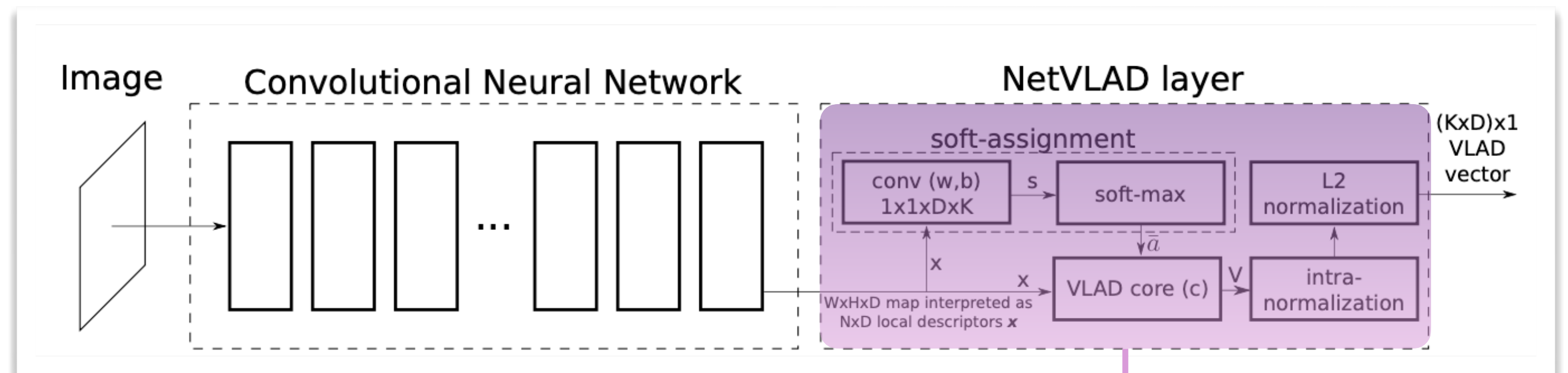
Network Architecture of PointNetVLAD.

Uy, Mikaela Angelina and Gim Hee Lee. "PointNetVLAD: Deep Point Cloud Based Retrieval for Large-Scale Place Recognition." *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2018): 4470-4479.

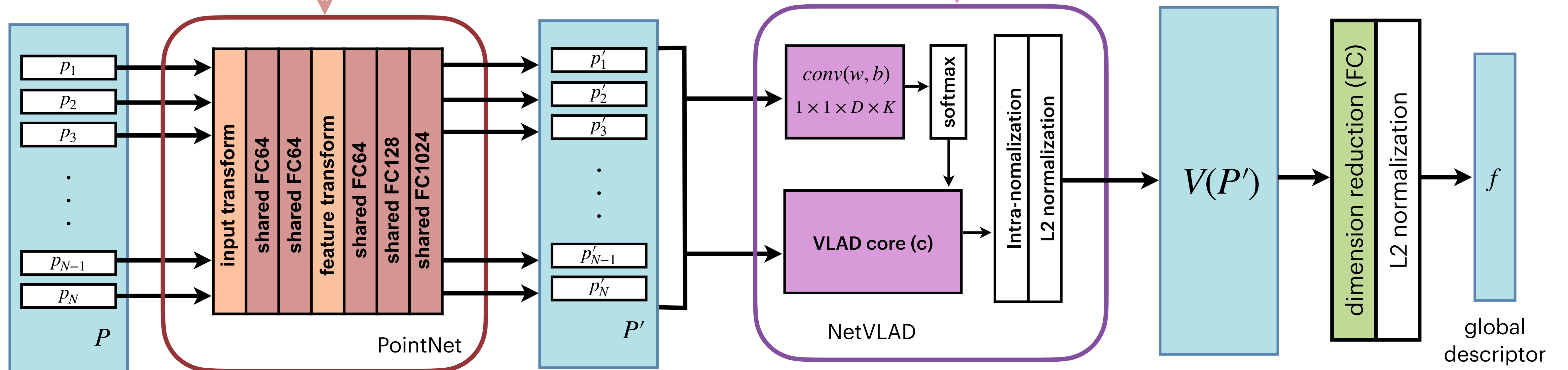
PointNetVLAD Backbone



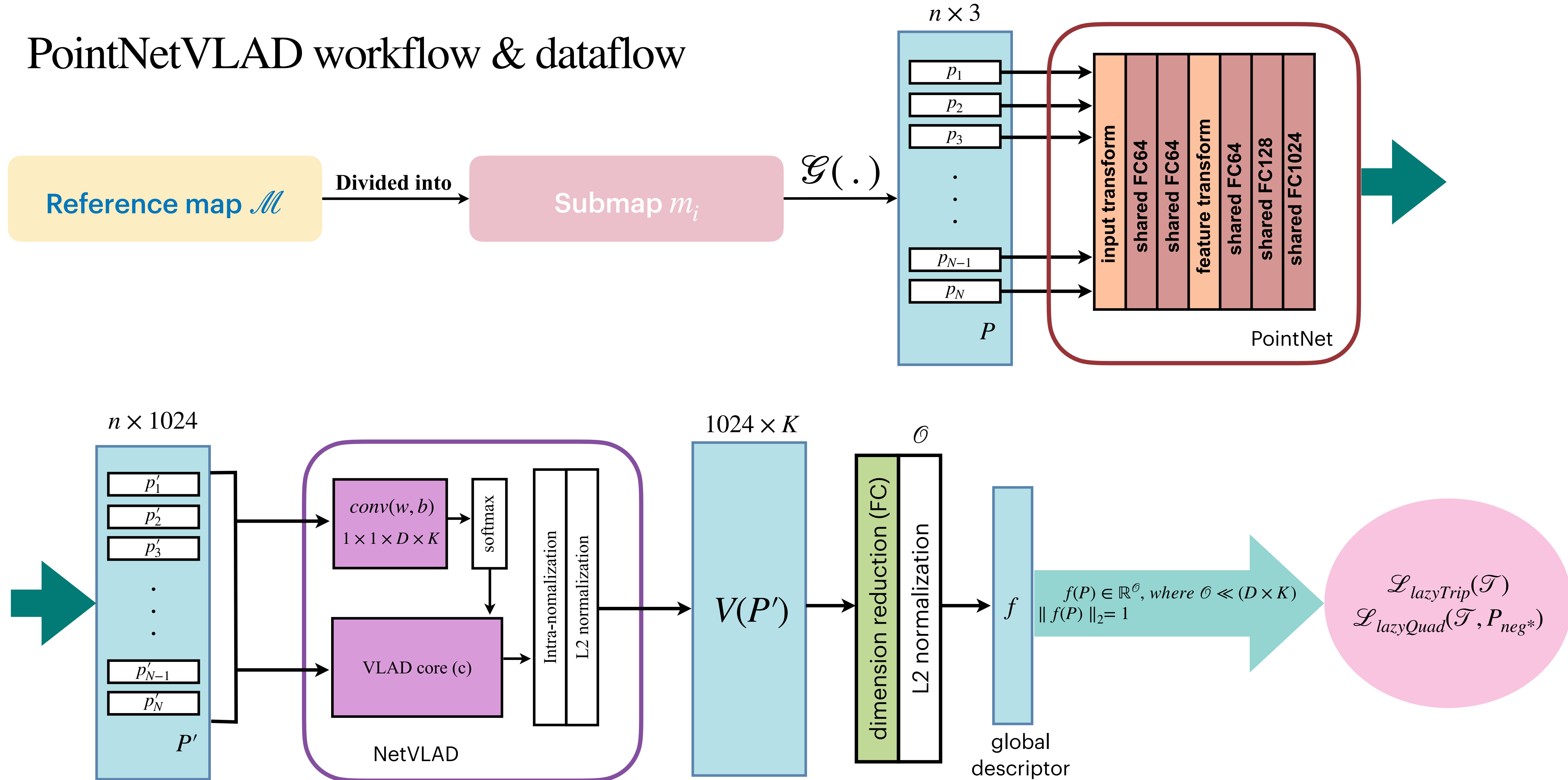
Network architecture of PointNet



CNN architecture with the NetVLAD layer

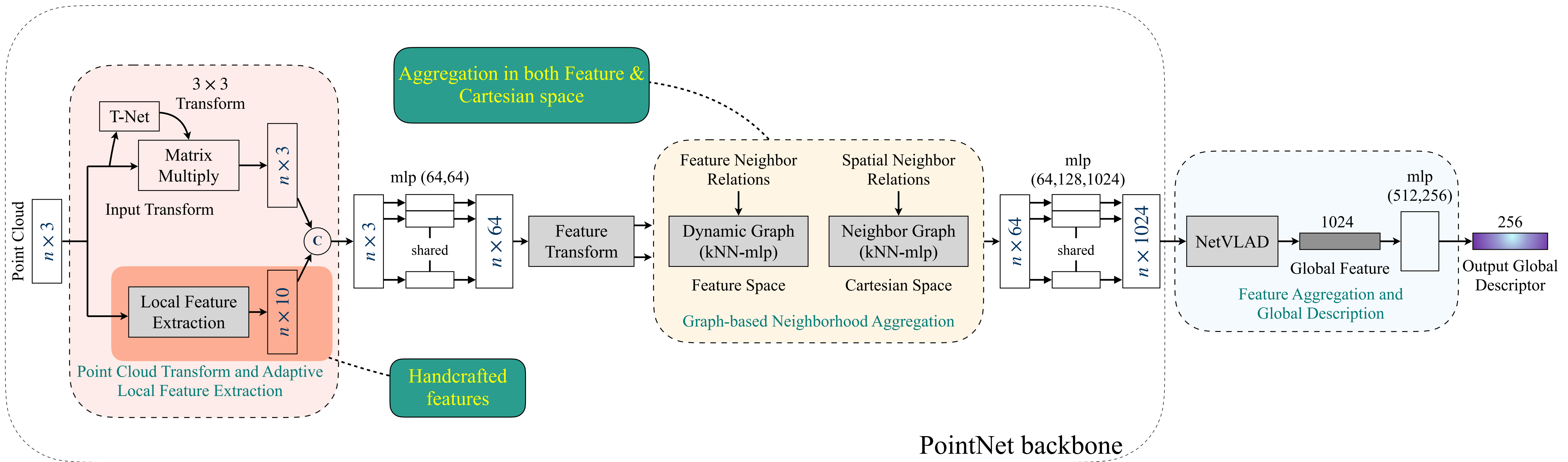


PointNetVLAD workflow & dataflow



LPD-Net

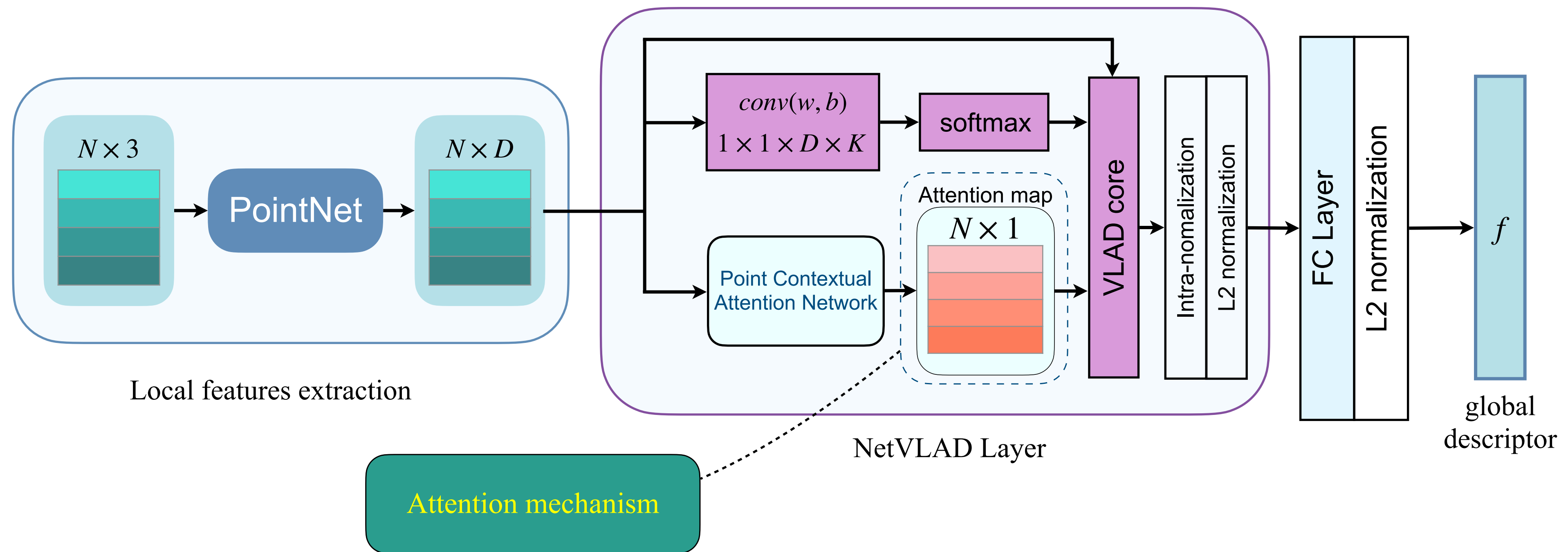
Handcrafted features, coordinate & feature space



Liu, Zhe et al. "LPD-Net: 3D Point Cloud Learning for Large-Scale Place Recognition and Environment Analysis." *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019): 2831-2840.

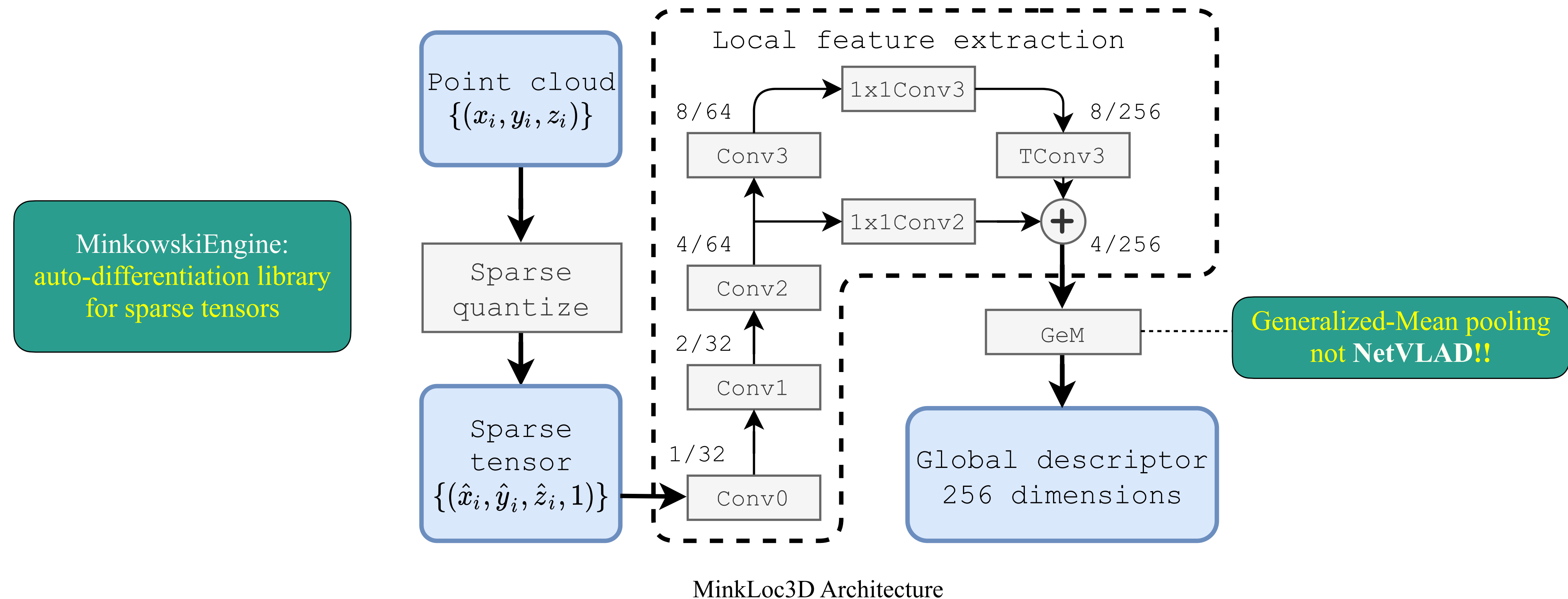
PCAN

The first one to introduce the attention mechanism



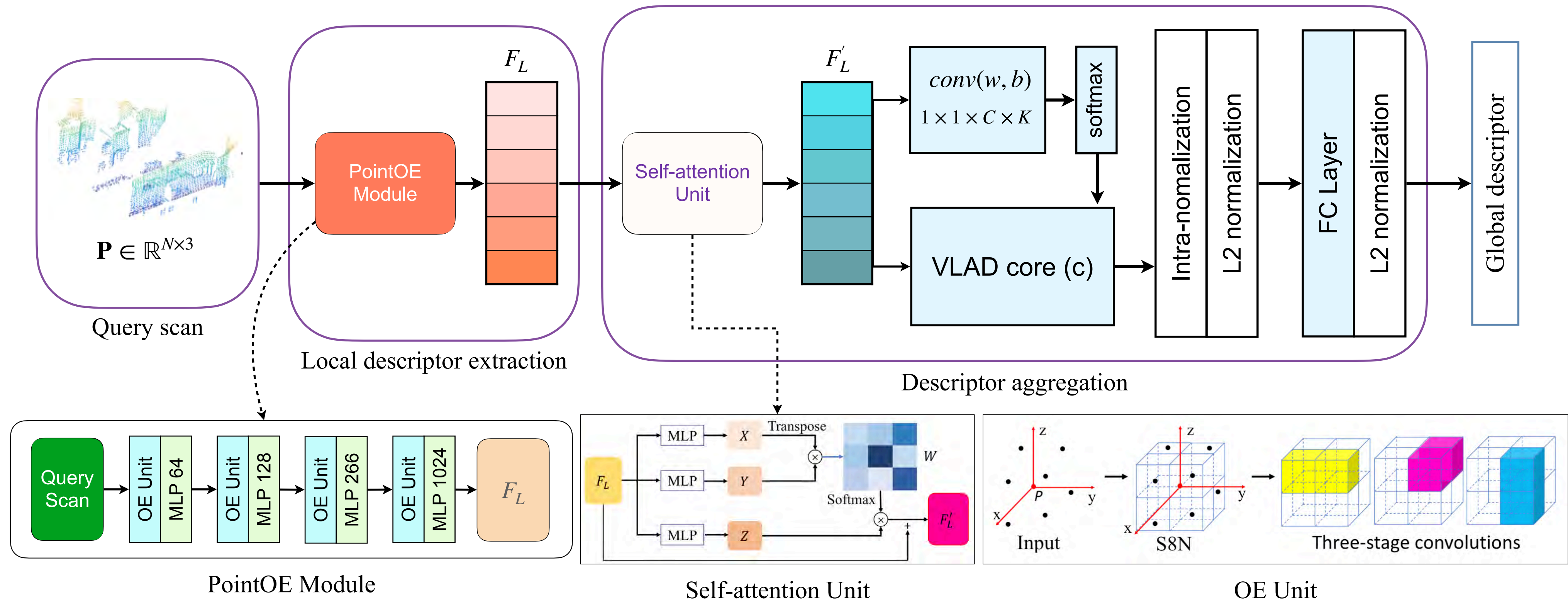
MinkLoc3D

A **novel** backbone for 3D VPR



SOE-Net

Introduce self-attention & a new loss function



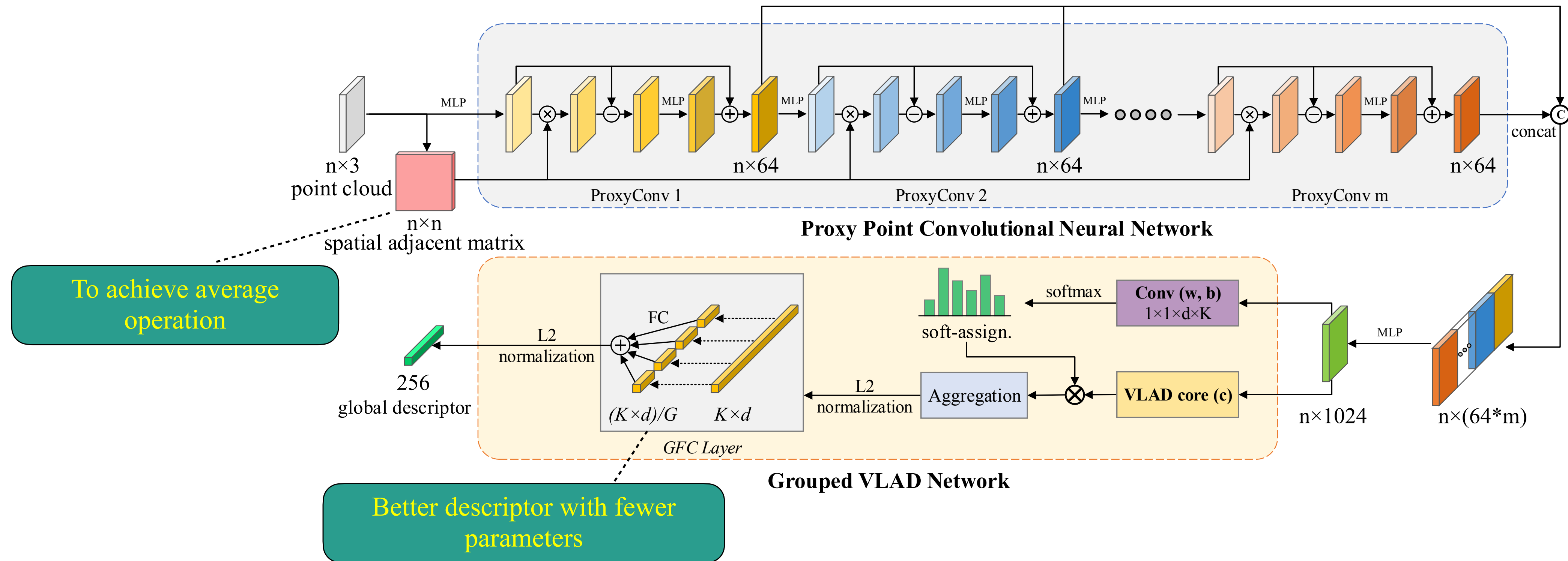
A New Loss Function for VPR:

Hardest **P**ositive **H**ardest **N**egative quadruplet loss (**HPHN loss**)

$$L_{HPHN} = \left[\left\| f(\delta_a) - f(\delta_{hp}) \right\|_2^2 - d_{hn} + \gamma \right]_+$$

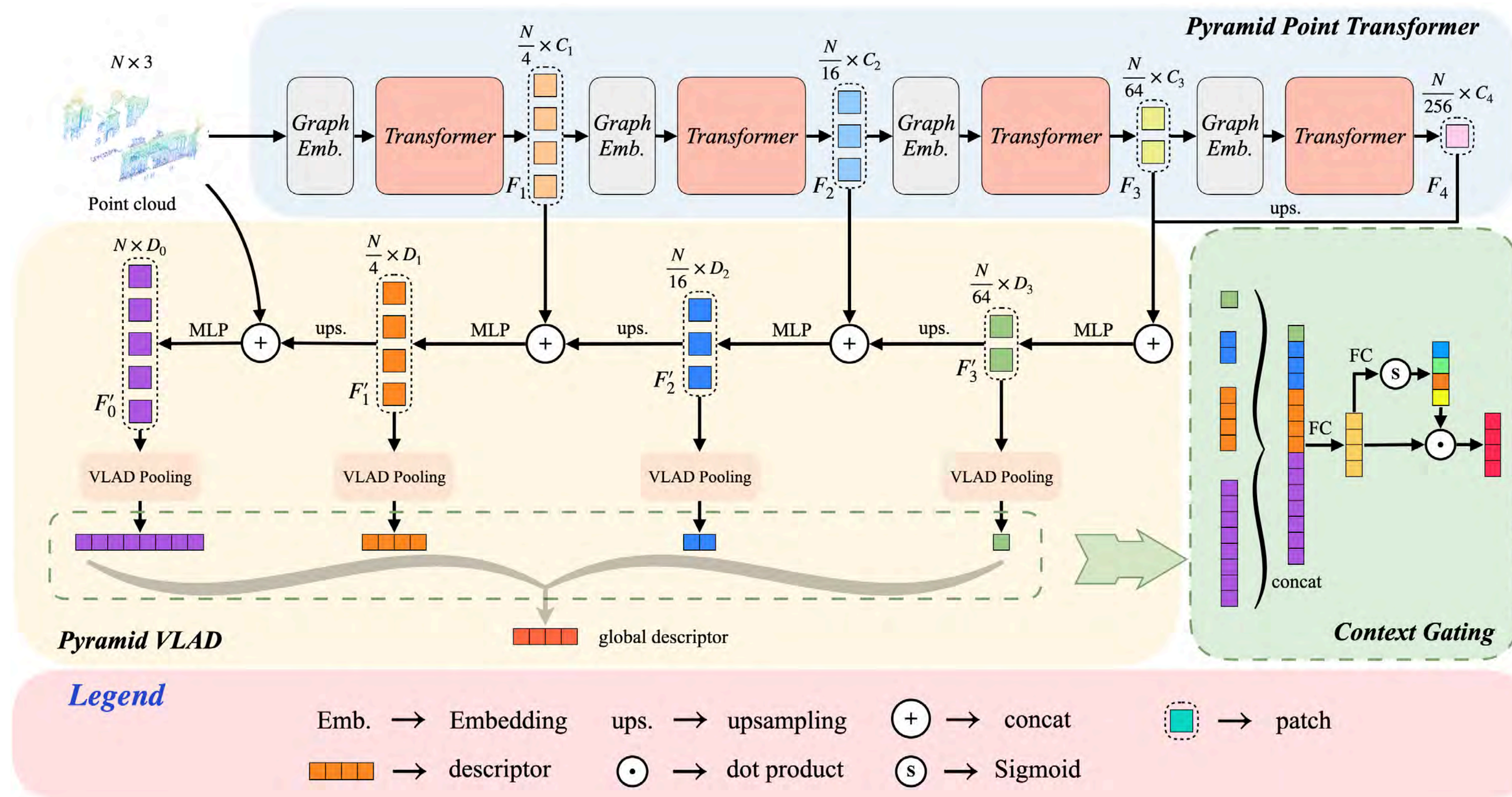
EPC-Net

ProxyConv & grouped VLAD



PPT-Net

A master of various algorithms



The pipeline of the pyramid point cloud transformer network (PPT-Net)

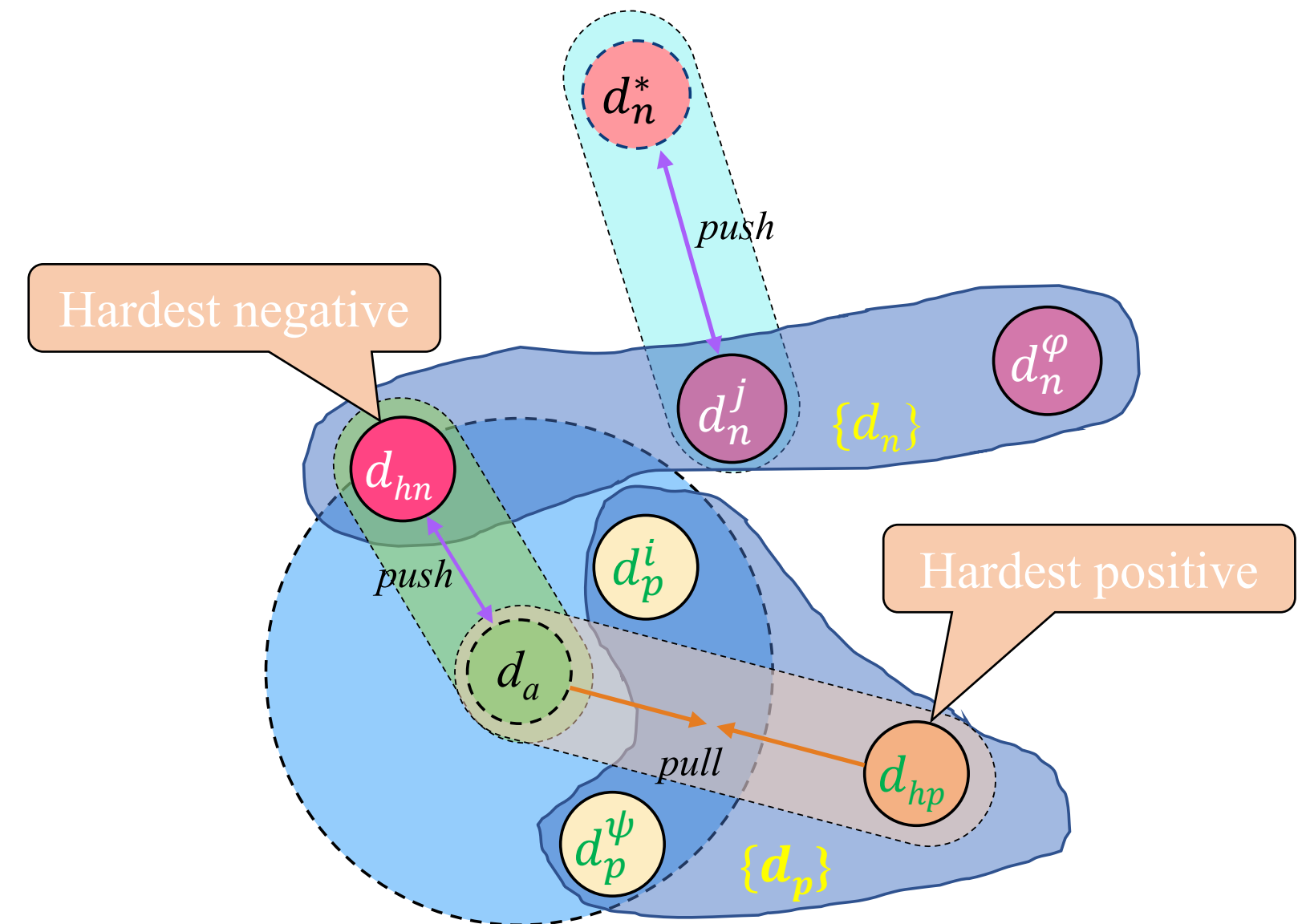
Hui, Le et al. "Pyramid Point Cloud Transformer for Large-Scale Place Recognition." *2021 IEEE/CVF International Conference on Computer Vision (ICCV)* (2021): 6078-6087.

IV. Methodology

HPHN quadruplet loss

Hardest **P**ositive **H**ardest **N**egative

- For a lazy quadruplet $Q_l = (d_a, \{d_p\}, \{d_n\}, d_n^*)$,
where d_a is the anchor point cloud,
 $\{d_p\}$ is a collection of ψ positive point clouds,
 $\{d_n\}$ is a collection of φ negative point clouds,
 d_n^* is a randomly sampled point cloud, structurally dissimilar to d_a , d_p and d_n .



HPHN quadruplet loss

- The hardest positive point cloud d_{hp} is the **least structurally similar** to the anchor point cloud.
- The hardest negative point cloud is the **most structurally similar** to the anchor point cloud or the randomly sampled point cloud d_n^* .

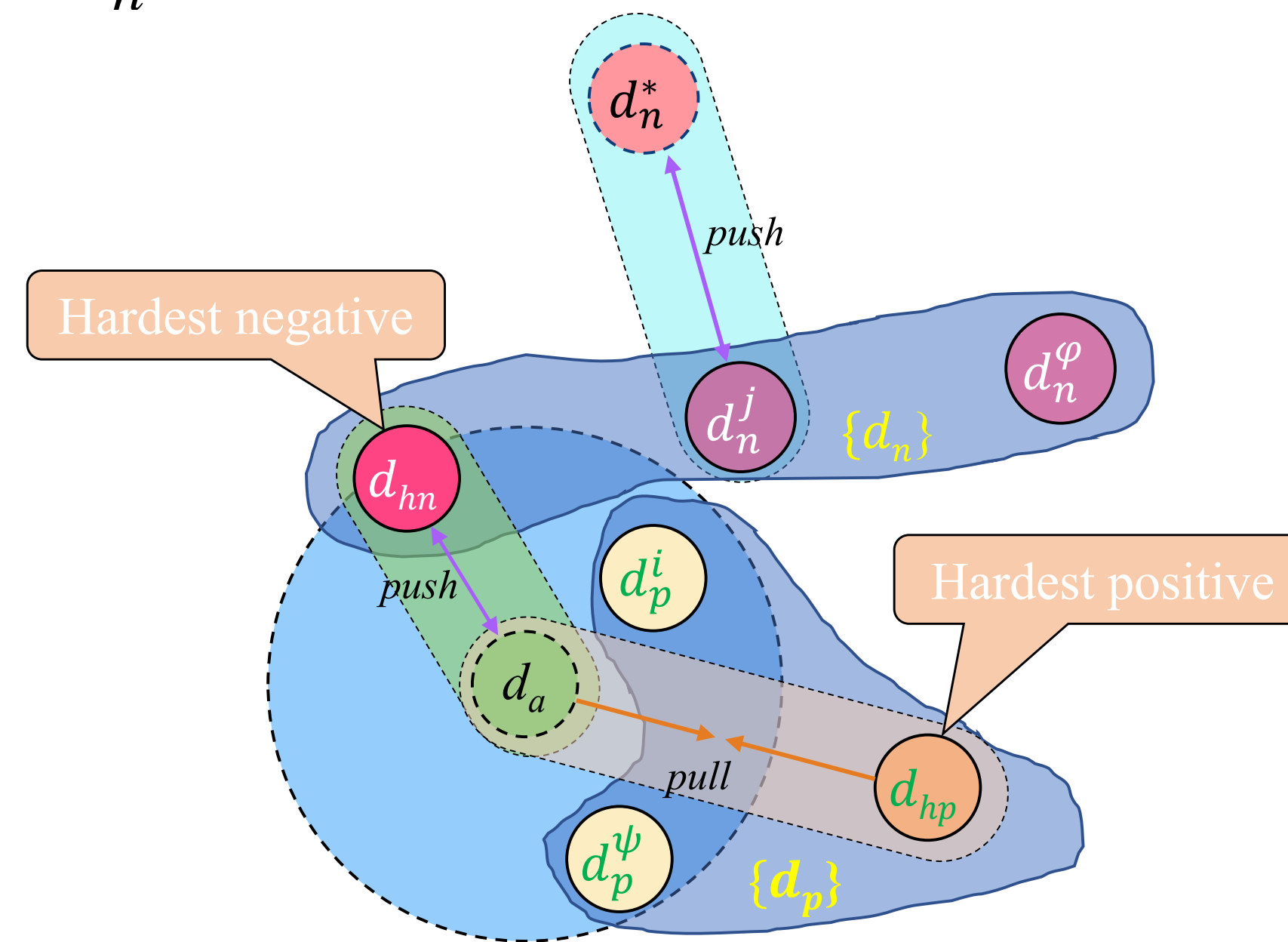
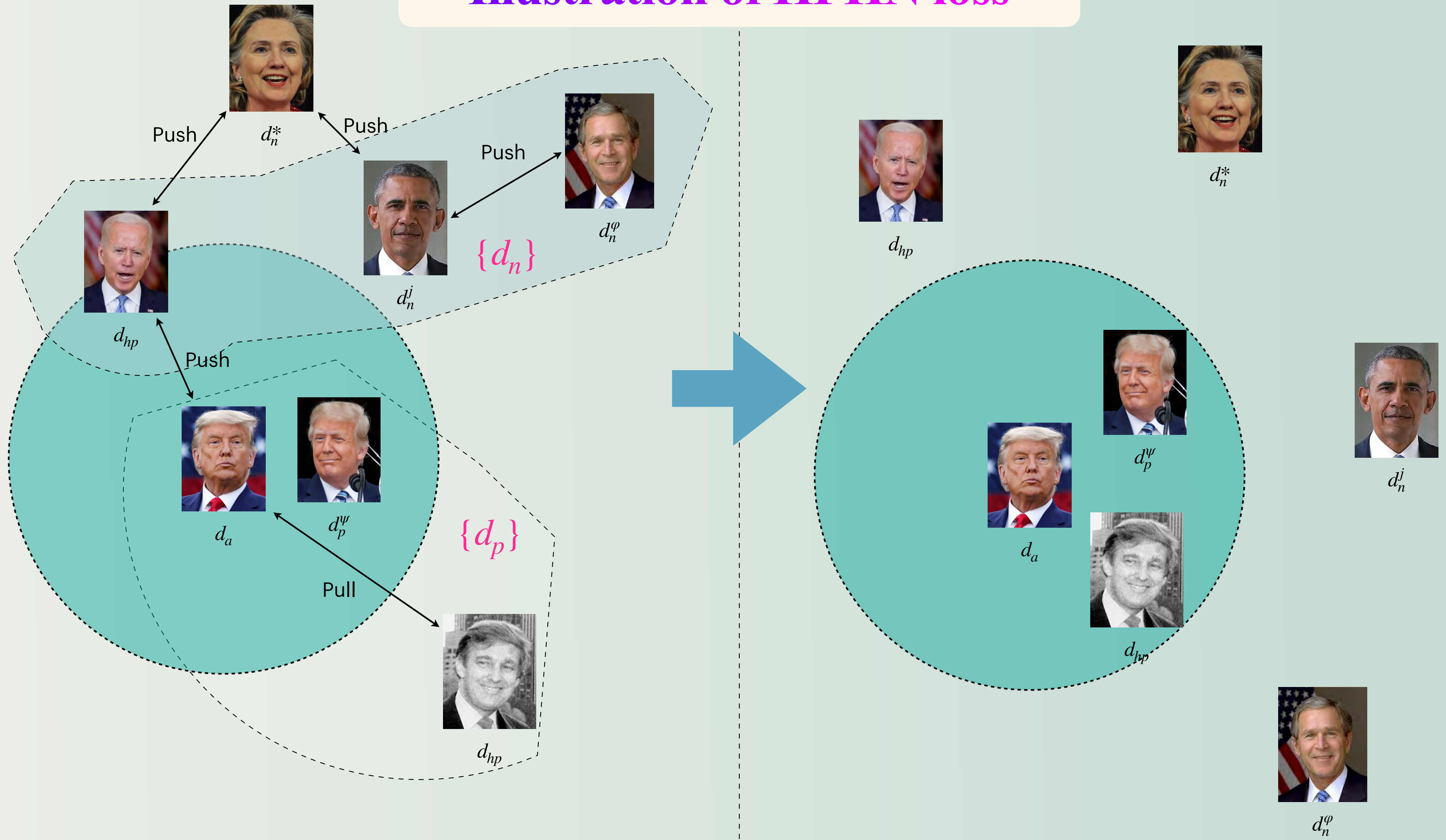


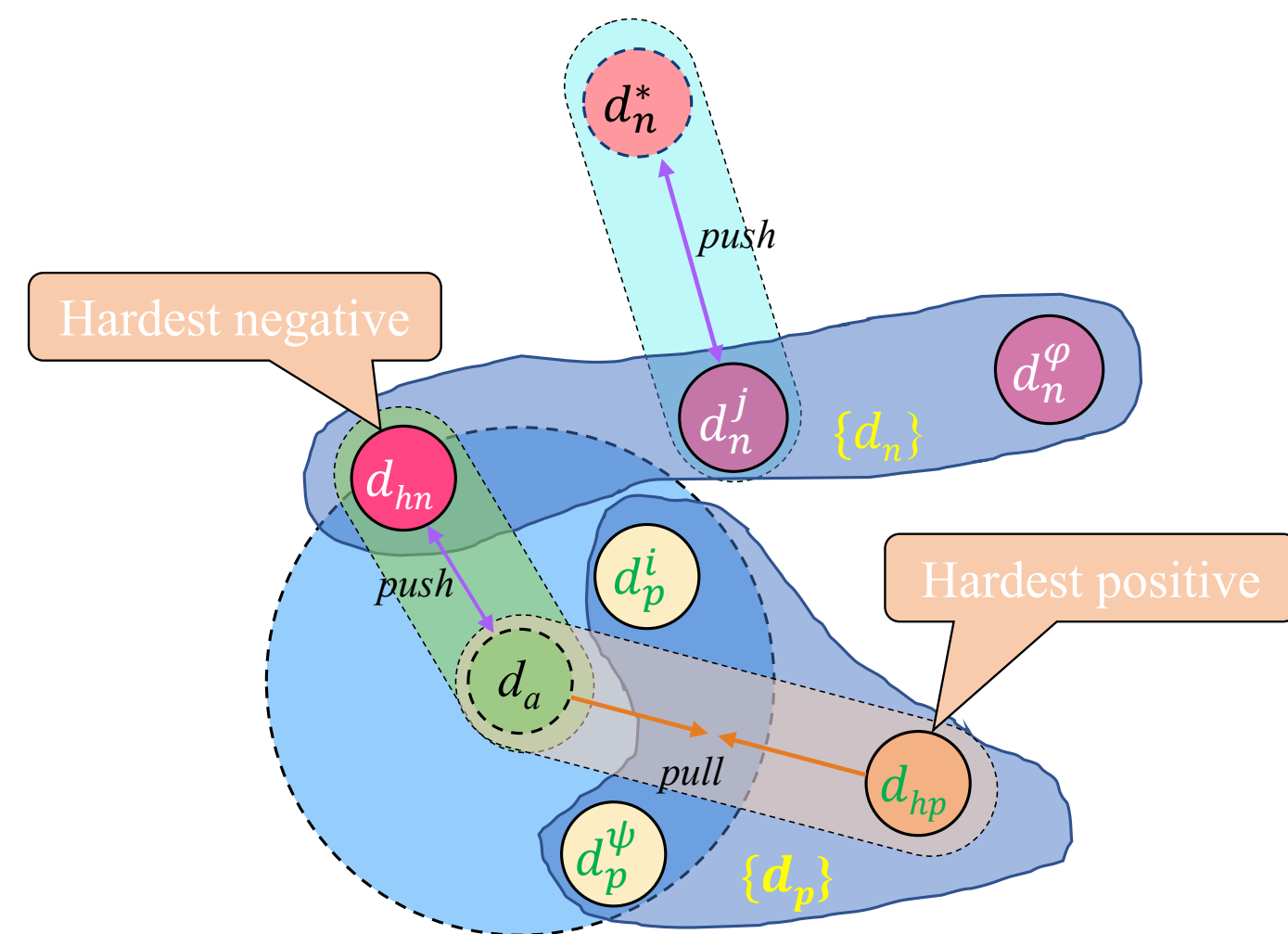
Illustration of HPHN loss



HPHN quadruplet loss

- In conclusion, the final HPHN quadruplet loss can be formulated as:

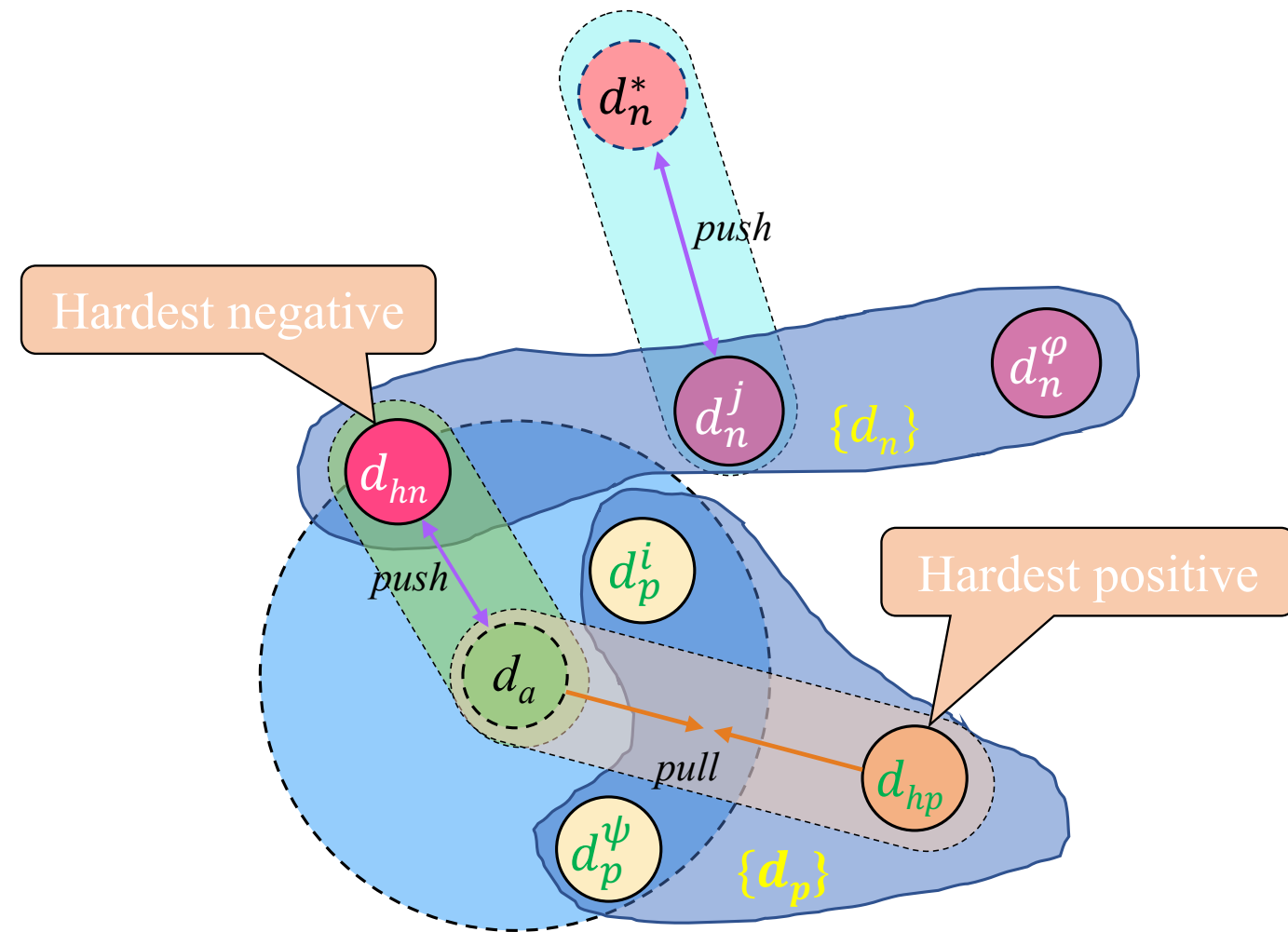
$$L_{HPHN} = \left[\left\| f(d_a) - f(d_{hp}) \right\|_2^2 - D_{hn} + \gamma \right]_+$$



$$= \left[\left(\text{pull } d_a, d_{hp} \right)^2 - \min \left\{ \left(\text{push } d_{hn}, d_a \right)^2, \left(\text{push } d_n^*, d_n^j \right)^2 \right\} + \gamma \right]_+$$

Scaled-HPHN quadruplet loss

Scaled-HPHN loss



$$L_{HPHN} = \left[\left\| f(d_a) - f(d_{hp}) \right\|_2^2 - D_{hn} + \gamma \right]_+$$

$$= \left[\left(\text{pull } d_a, d_{hp} \right)^2 - \min \left\{ \left(\text{push } d_{hn}, d_a \right)^2, \left(\text{push } d_n^*, d_n^j \right)^2 \right\} + \gamma \right]_+,$$

introduce a **scale factor** κ , then we have

$$L_{S-HPHN} = \left[\left\| f(d_a) - f(d_{hp}) \right\|_2^2 - \kappa D_{hn} + \gamma \right]_+$$

V. Experiments

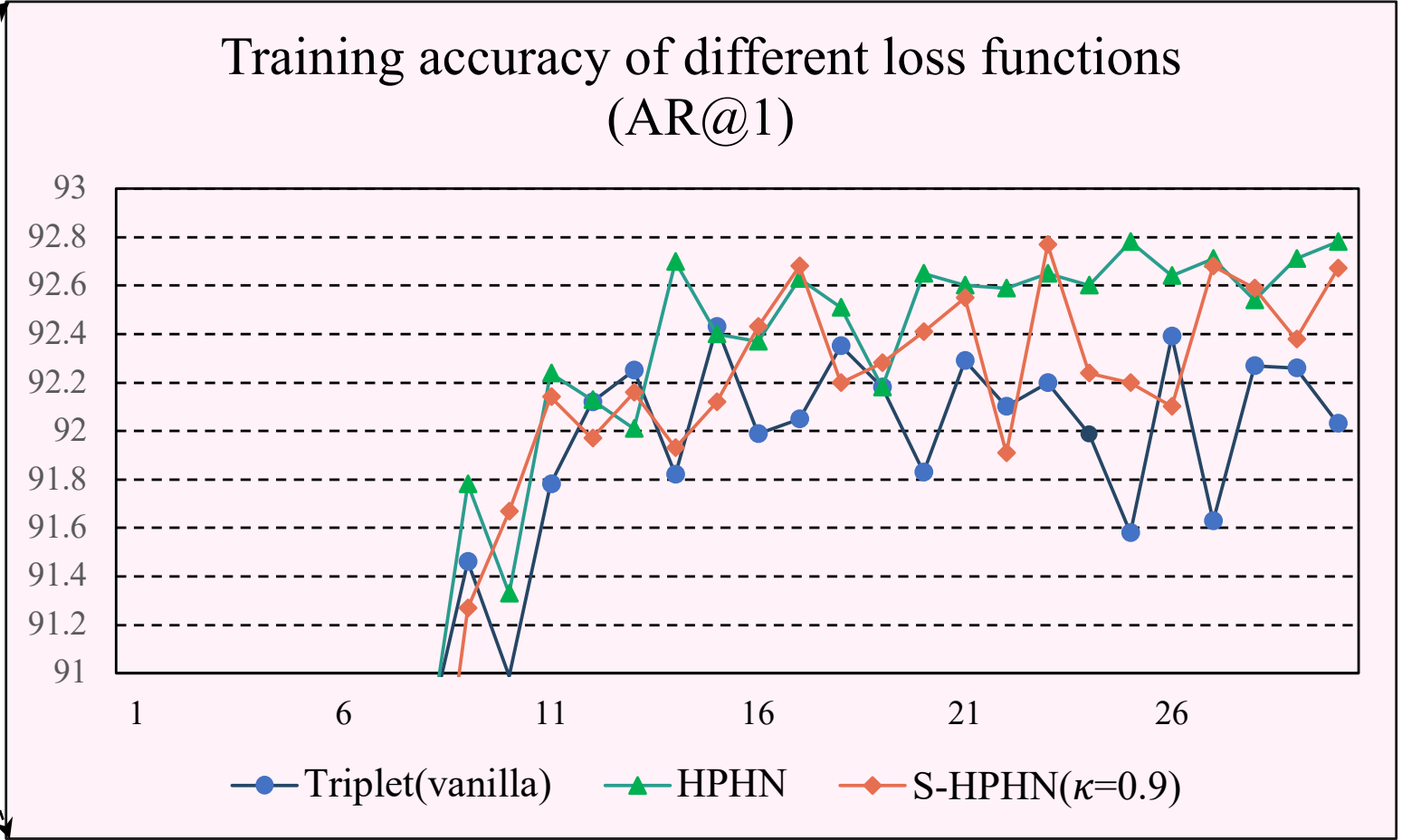
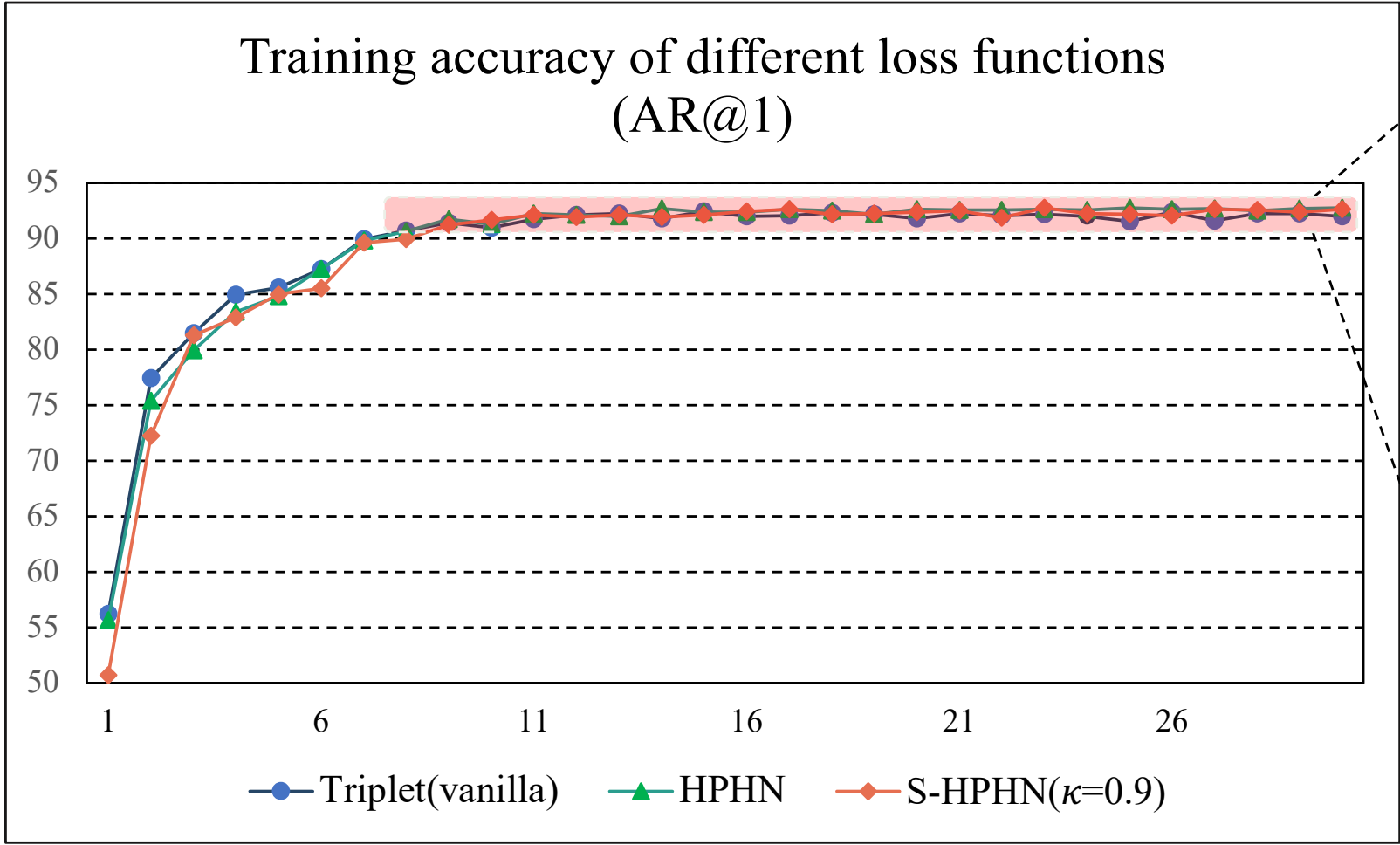
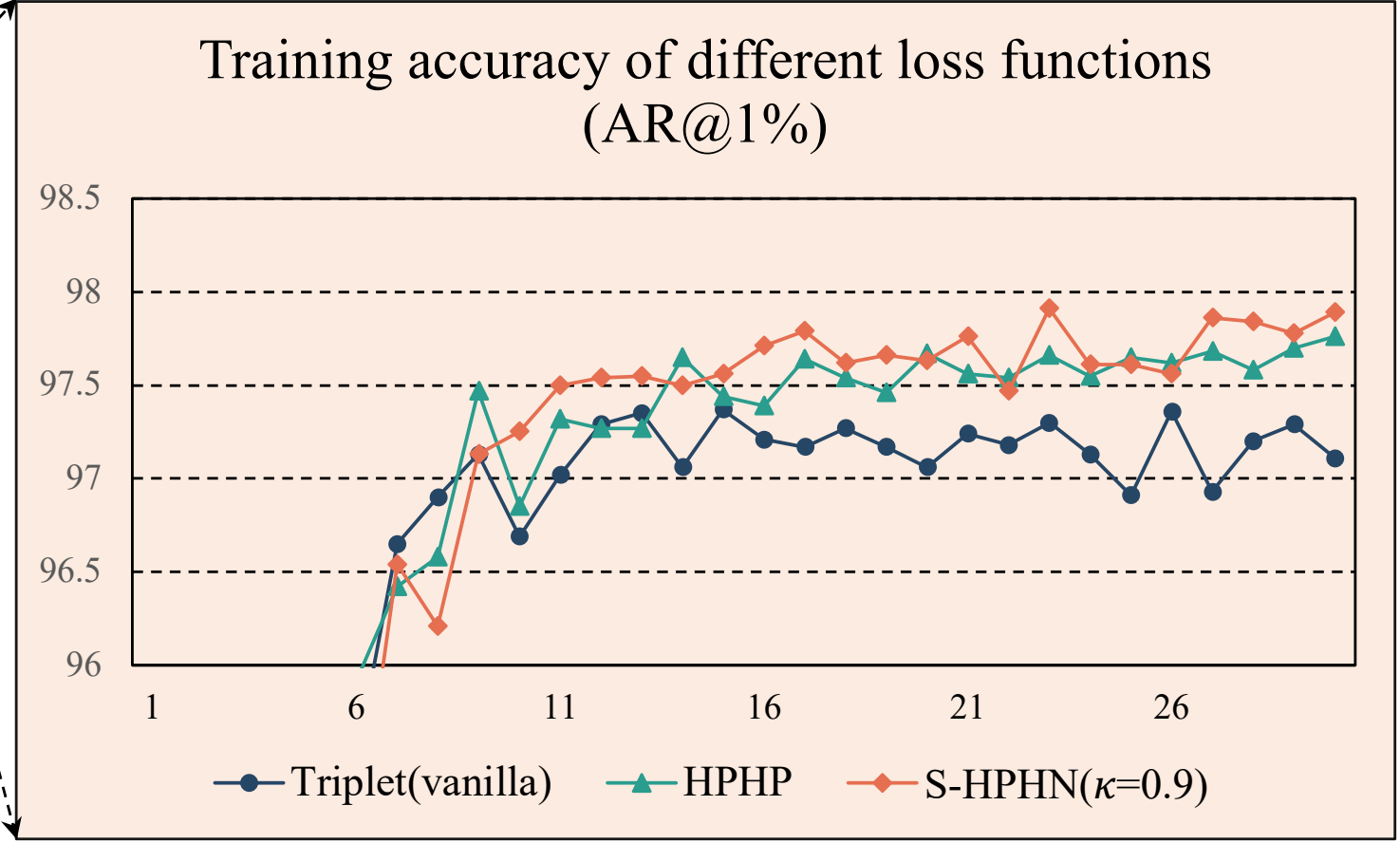
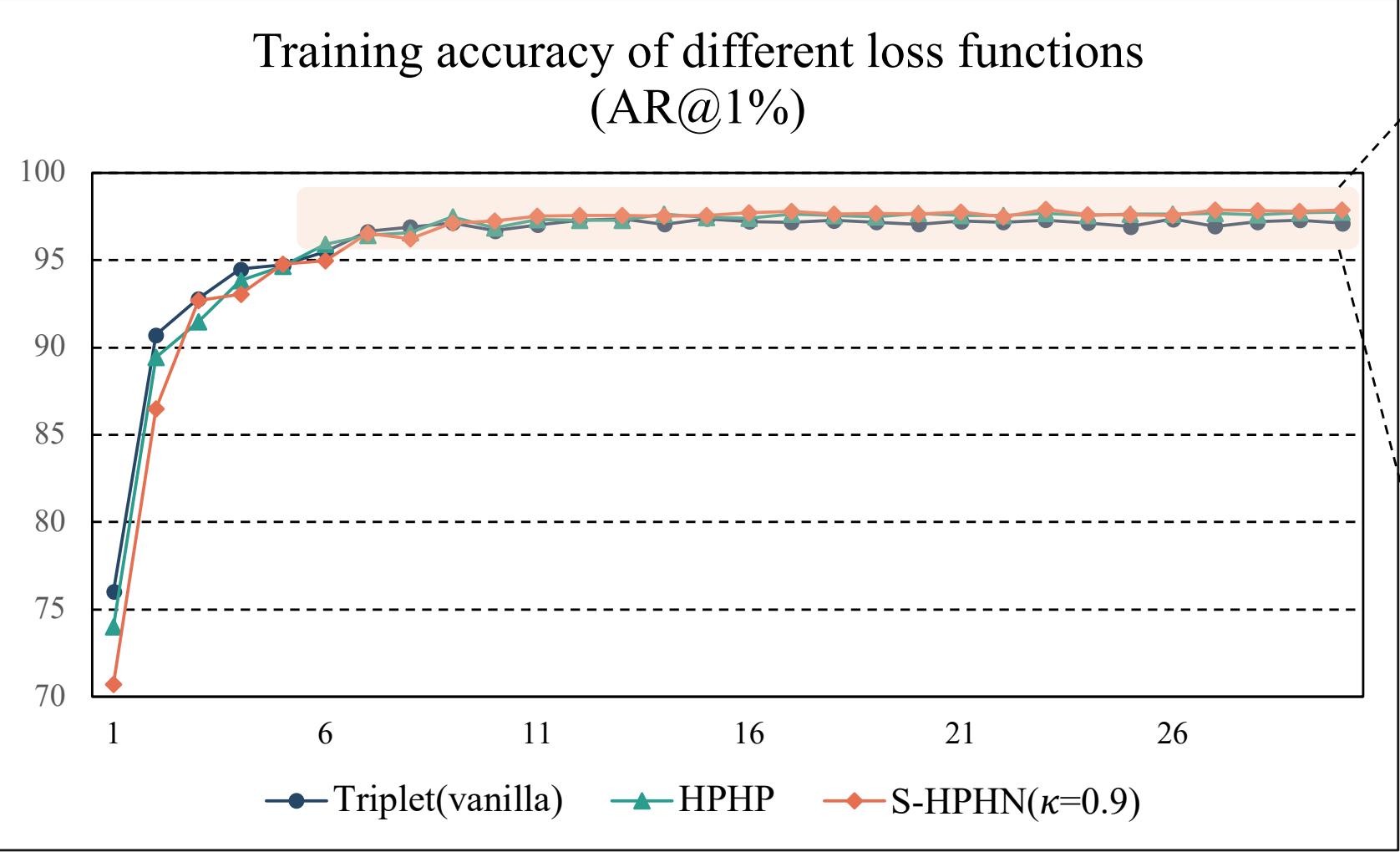
HPHN in PPT-Net

Introduce HPHN to PPT-Net

Loss	AR@1%	AR@1
Triplet (vanilla)	97.11	92.03
HPHN	97.76 (+0.65)	92.78 (+0.75)

Ablation study

κ	AR@1%	AR@1
0.9	97.89	92.67
1	97.76	92.78
1.001	97.60	92.65
1.01	96.98	91.59
1.1	97.17	91.99



Comparison between three losses

Loss	AR@1%	AR@1
Triplet (vanilla)	97.11	92.03
HPHN	97.76 (+0.65)	92.78 (+0.75)
S-HPHN ($\kappa=0.9$)	97.89 (+0.78)	92.67 (+0.64)

VI. Conclusion

- The HPHN loss is implemented in PPT-Net.
- Based on HPHN loss, a scale factor is introduced to propose the scaled-HPHN loss function.
- Experiments show that both HPHN and scaled-HPHN are better than the original triplet loss.

VII. Future work

Trainable scale factor

$$L_{S-HPHN} = \left[\left\| f(d_a) - f(d_{hp}) \right\|_2^2 - \kappa D_{hn} + \gamma \right]_+$$

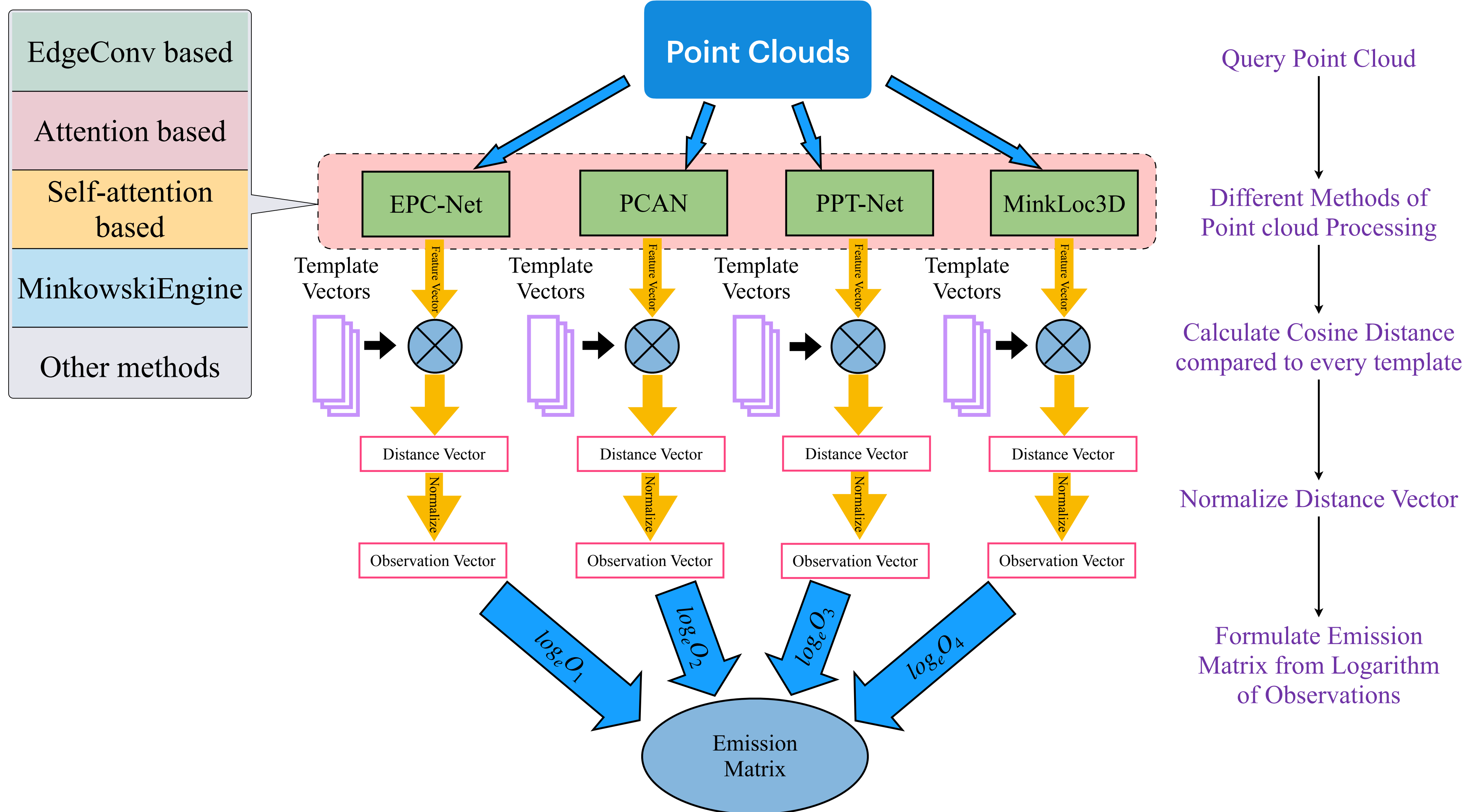
Make it trainable

Classification of 3D VPR Algorithms

EdgeConv based	Attention based	Self-attention based	MinkowskiEngine
LPD-Net EPC-Net PPT-Net	PCAN DH3D	SOE-Net PPT-Net	MinkLoc3D

Method	Parameters	FLOPs	Runtime per frame
PN_VLAD	19.78M	4.21G	20ms
PCAN	20.42M	7.73G	58ms
LPD-Net	19.81M	7.80G	28ms
MinkLoc3D	1.10M	1.81G	17ms
EPC-Net	4.70M	3.25G	20ms
PPT-Net	13.12M	3.23G	18ms

Multi-process Fusion



2 & A