



# BEV-Guided Multi-Modality Fusion for Driving Perception

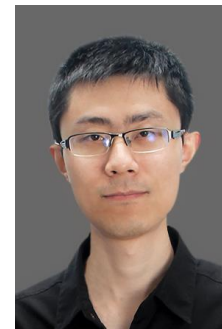
Yunze Man



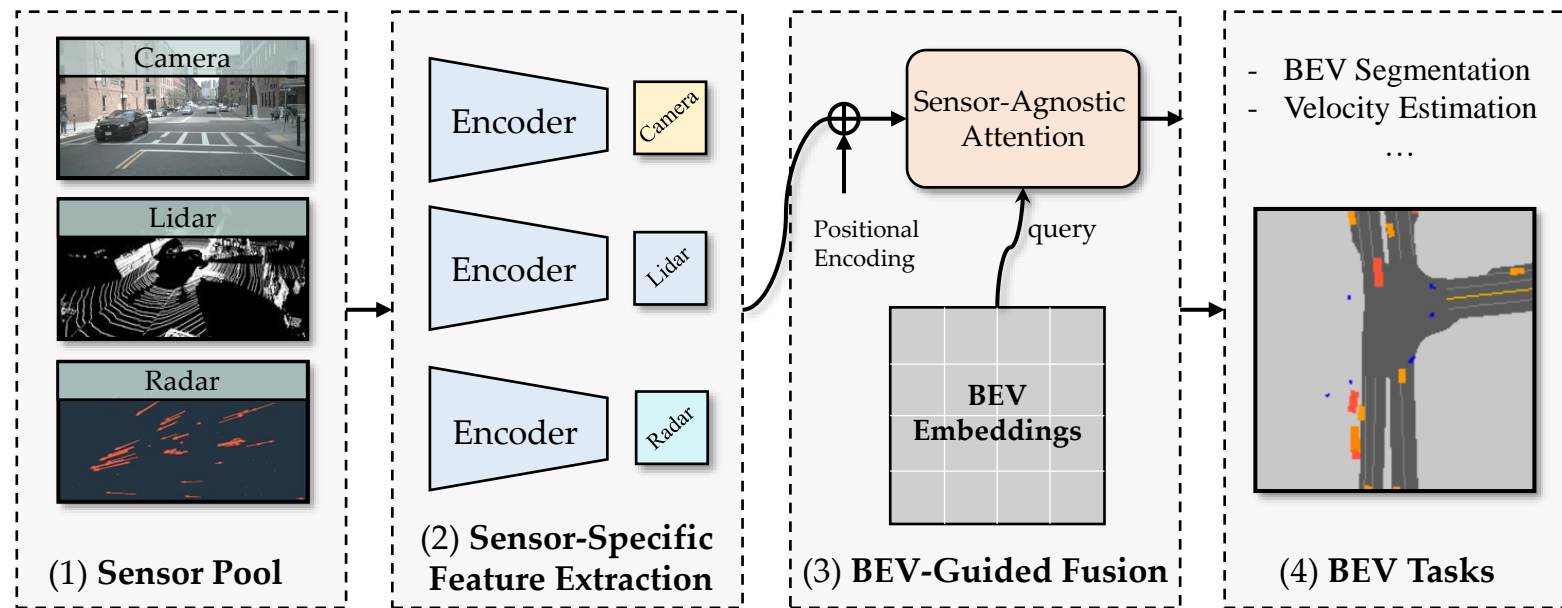
Liang-Yan Gui



Yu-Xiong Wang



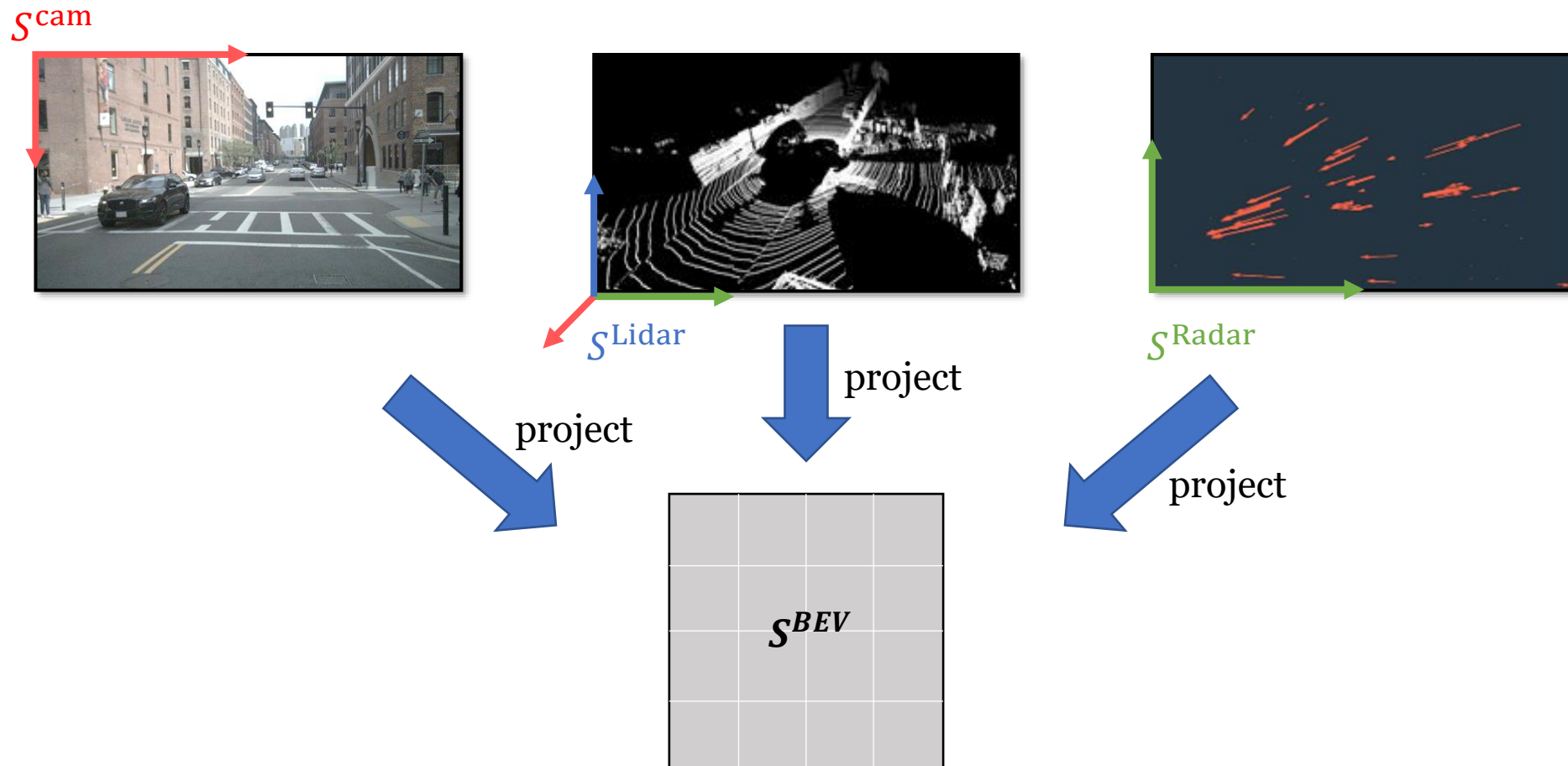
# Overview



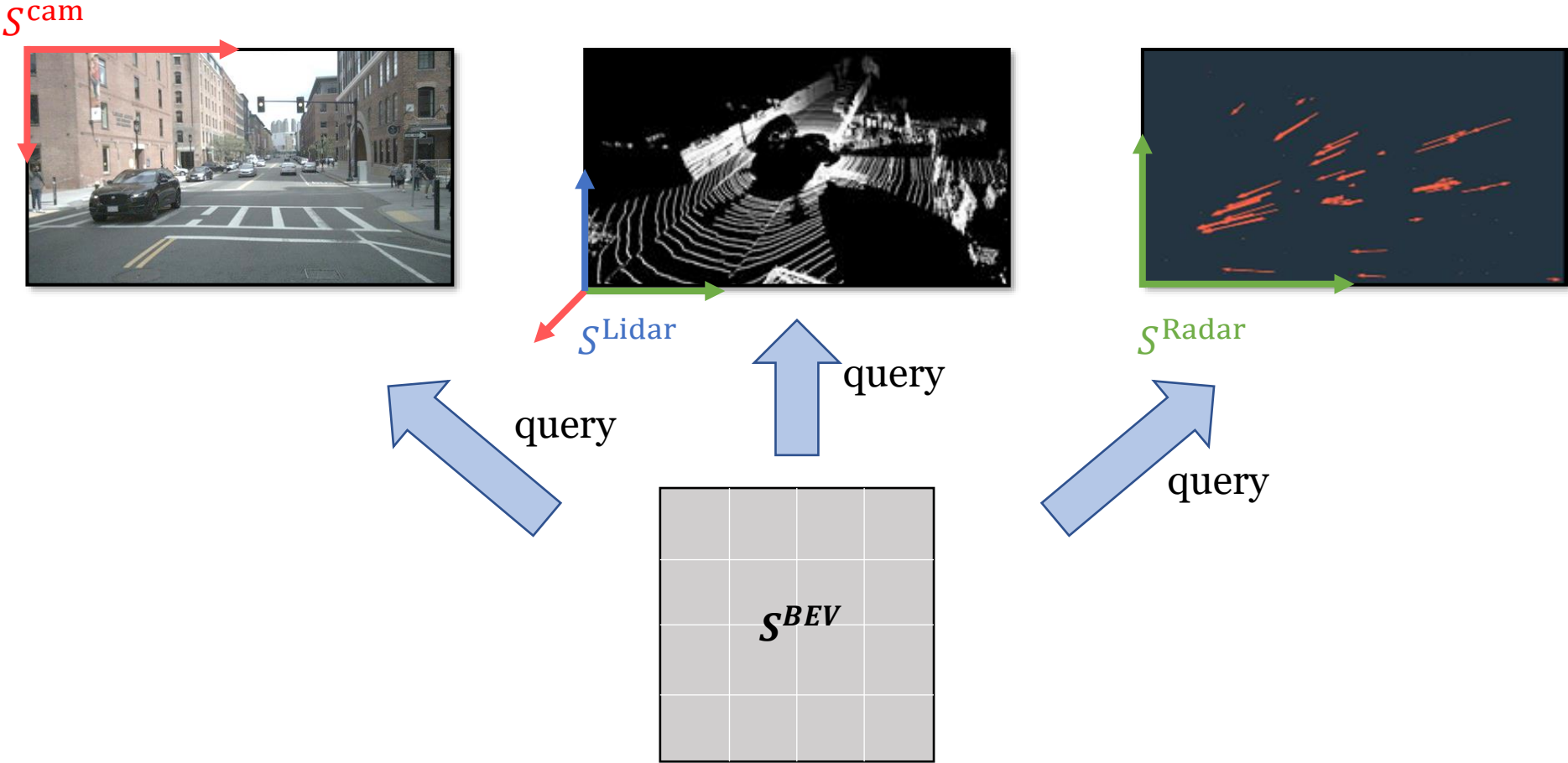
## BEVGuide

**Comprehensive** and **versatile** multi-modality fusion architecture for driving tasks

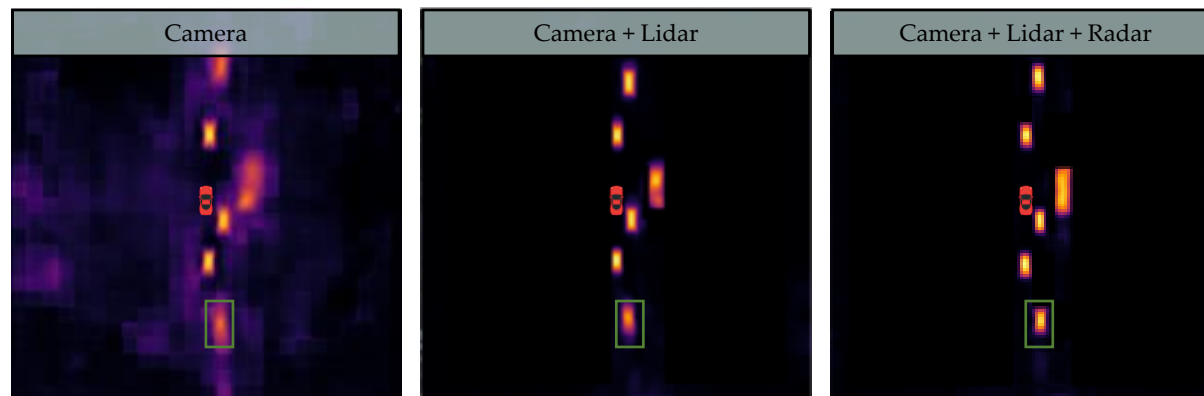
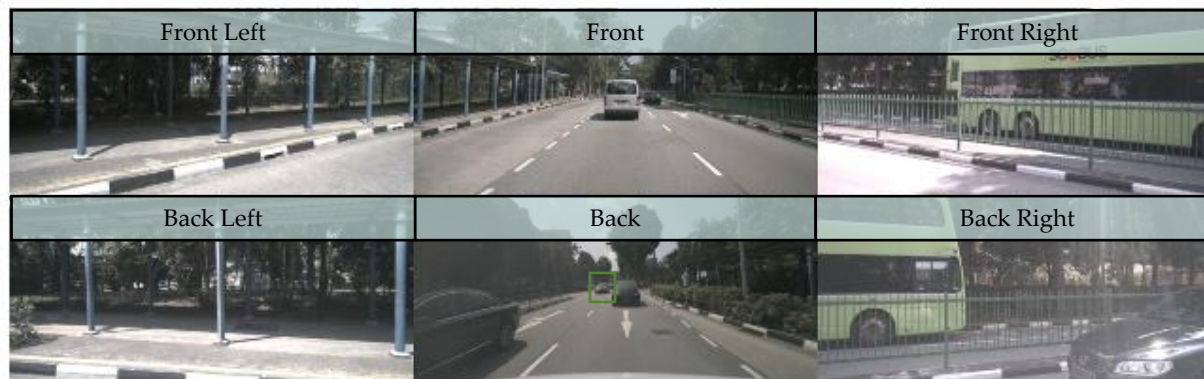
# Previous Fusion Methods: Project Sensors **into** the Unified Space



# Our Key Insight: Query Sensors from the Unified Space



# Superior Performance on Segmentation



# Superior Performance on Diverse Driving Tasks: Segmentation, Detection, and Velocity Estimation

<i>Segmentation</i>	C R L	Vehicles ↑	Roads ↑
Cross-view	✓	36.0	74.3
FUTR3D	✓ ✓	46.6	-
Simple-BEV	✓ ✓	60.8	-
BEVFusion	✓ ✓	-	85.5
X-Align	✓ ✓	-	86.8
<b>BEVGuide</b>	✓ ✓ ✓	<b>79.0</b>	<b>86.9</b>

<i>Detection</i>	C R L	mAP ↑	NDS ↑
FUTR3D	✓ ✓	35.0	45.9
BEVGuide*	✓ ✓	42.1	53.7
BEVFusion	✓ ✓	68.5	71.4
BEVGuide*	✓ ✓	68.9	71.4
<b>BEVGuide</b>	✓ ✓ ✓	<b>69.3</b>	<b>71.5</b>

<i>Velo. Estimation</i>	C R L	P-AVE ↓
Cross-view	✓	2.13
PointPainting	✓ ✓	1.90
BEVGuide*	✓ ✓	1.63
<b>BEVGuide</b>	✓ ✓ ✓	<b>0.81</b>

# Summary

- BEVGuide, a **comprehensive** and **versatile** multi-modality fusion architecture
- Easily adapt to different **sensor combinations**
- Achieved **state-of-the-art** performance on various driving tasks

# Check Our Project Website



<https://yunzeman.github.io/BEVGuide/>



# Multi-Modality is Everywhere

Music videos: **Video + Audio**



Communication : **Video + Language + Audio + etc.**



Reading: **Language + Video + etc.**



# Multi-Modality is Everywhere

Music videos: **Video + Audio**



Communication : **Video + Language + Audio + etc.**

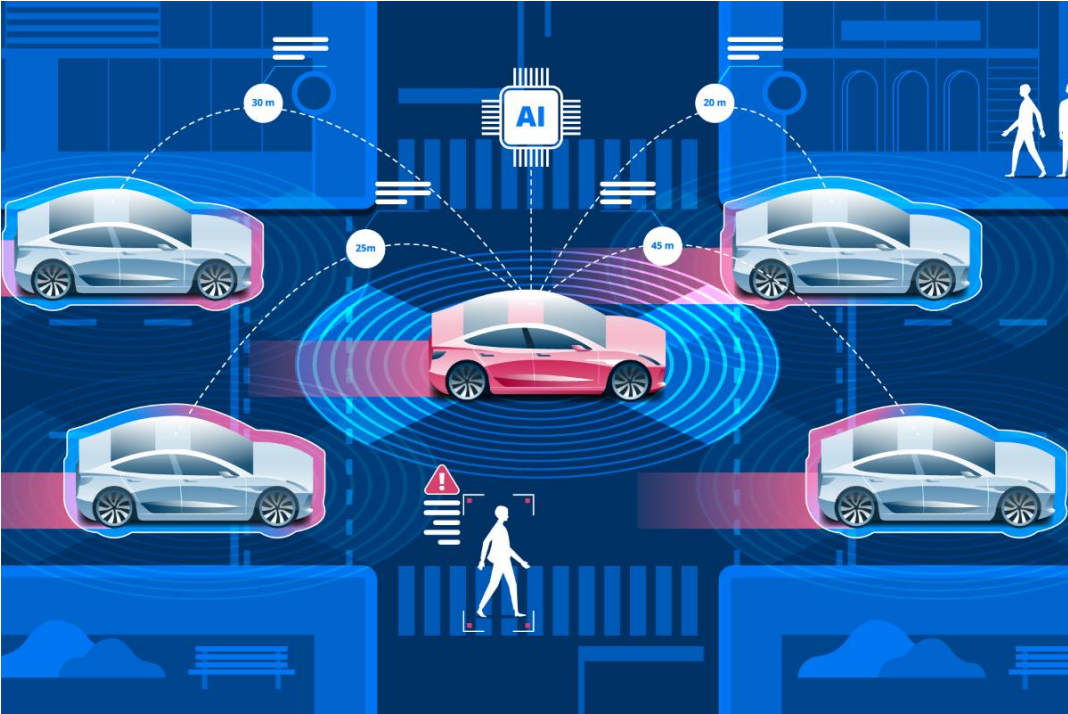


Reading: **Language + Video + etc.**



**Multi-Modality:**  
A more **complete** and **diverse** experience

# Multi-Modality in Driving Scenarios



Camera



Lidar



Radar

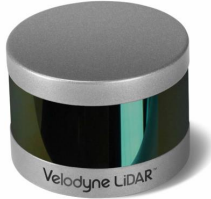
# Motivation of BEVGuide (Limitation of Prior Work)

- **Inflexible** design for different sensors (**Limited** and **fixed**)
- Do not **dynamically** change **weights** of different sensors based on diverse input sample
- Overlook **Radar** sensor and its unique properties

# Flexible Sensor Combination



+



+



+



+



+



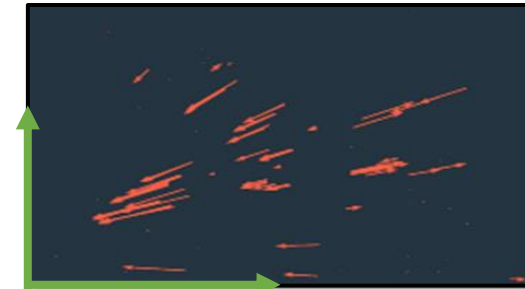
.....

# Different sensors are natively represented in different coordinates

$\zeta^{\text{cam}}$



$\zeta^{\text{Lidar}}$



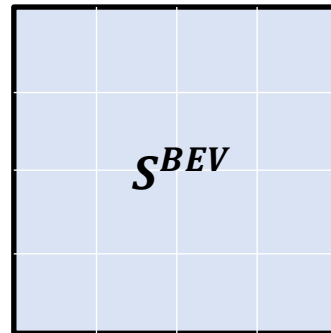
$\zeta^{\text{Radar}}$

# Bridge coordinate systems with BEV

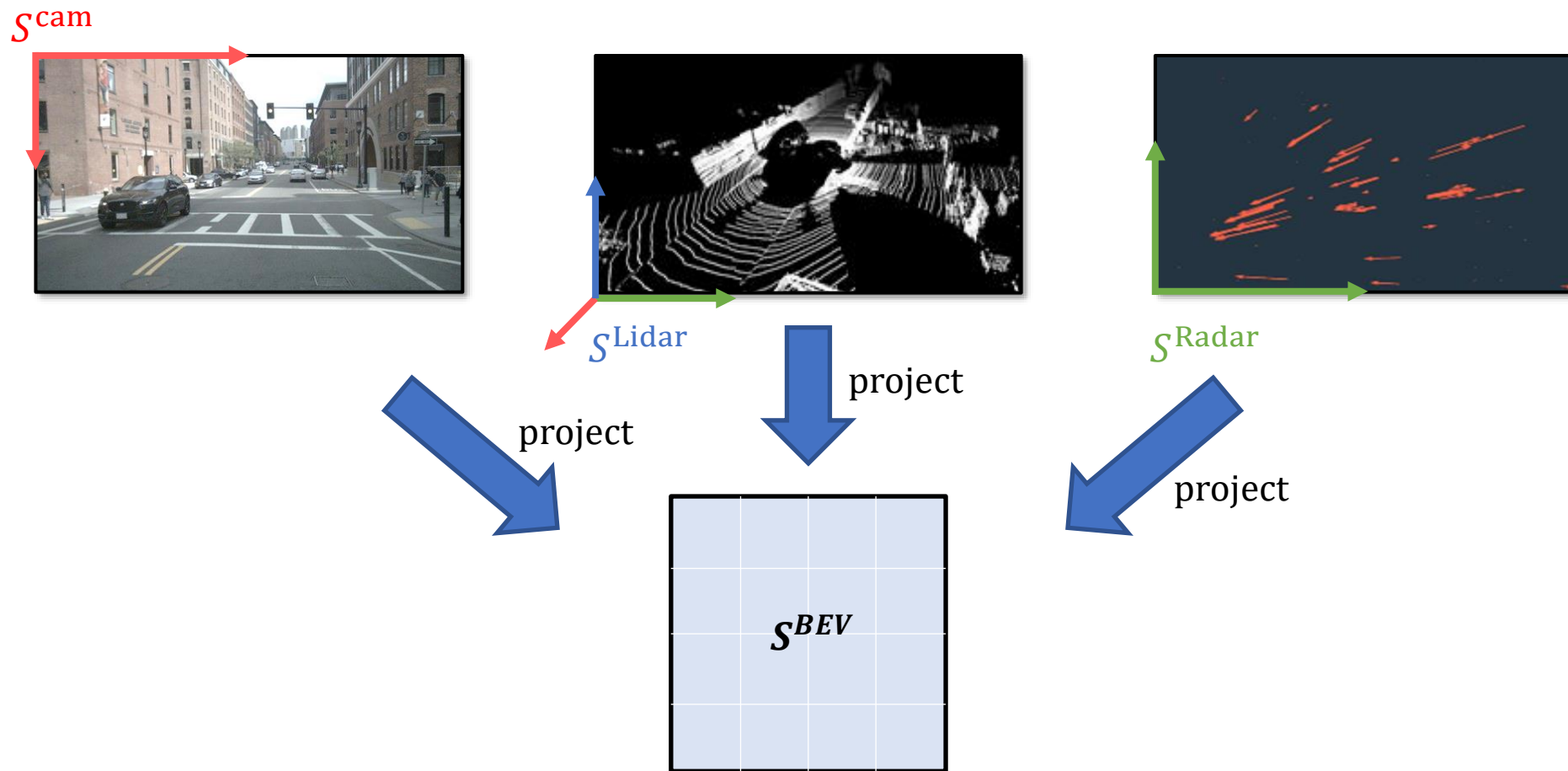
$\zeta_{cam}$



- Almost everything stands and moves on the **ground plane**
- **Simplify** 3D into 2D
- BEV is an **expressive** and **concise** coordinate system



# Previous Fusion Methods: Project Sensors into the Unified Space

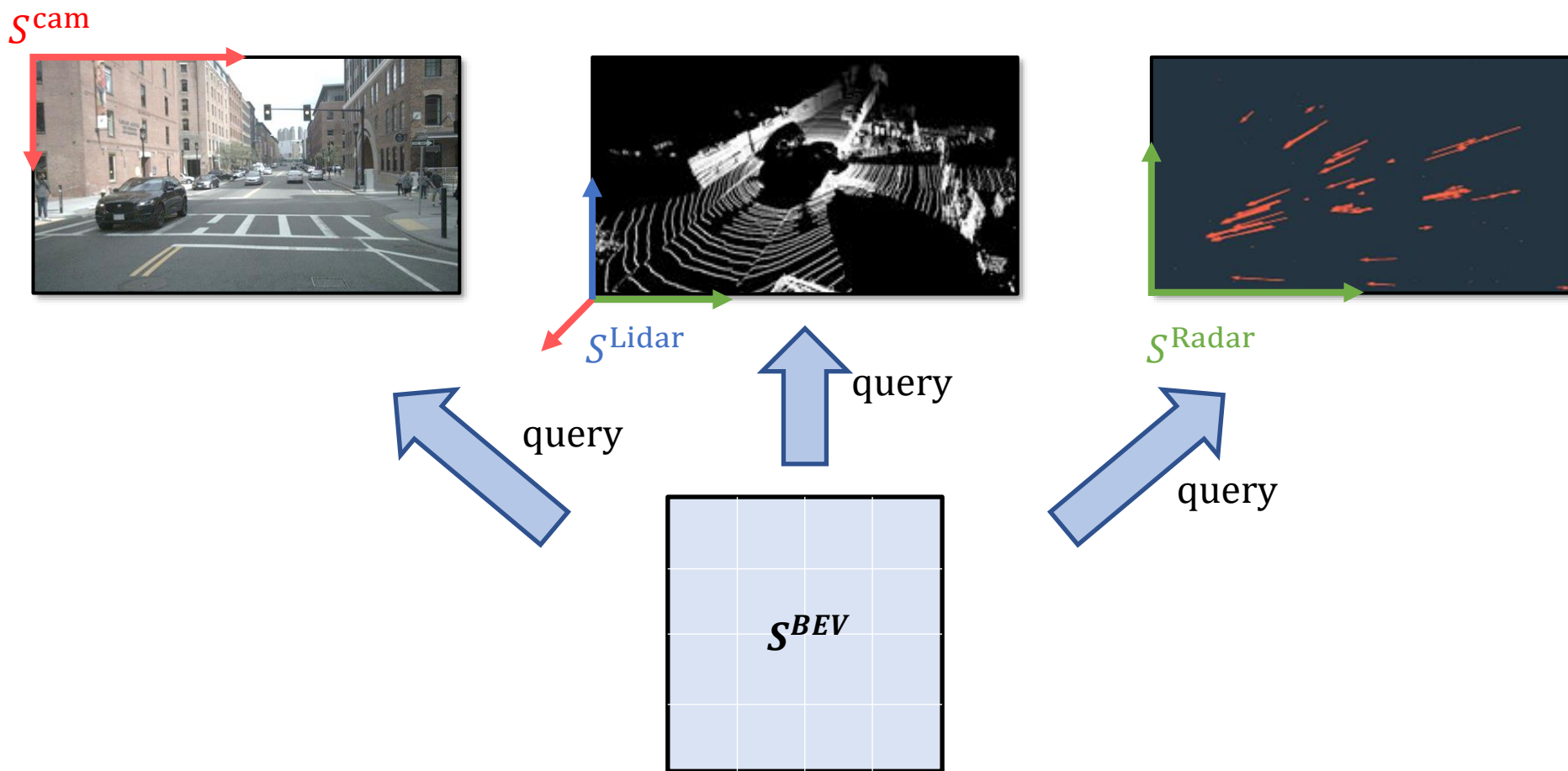




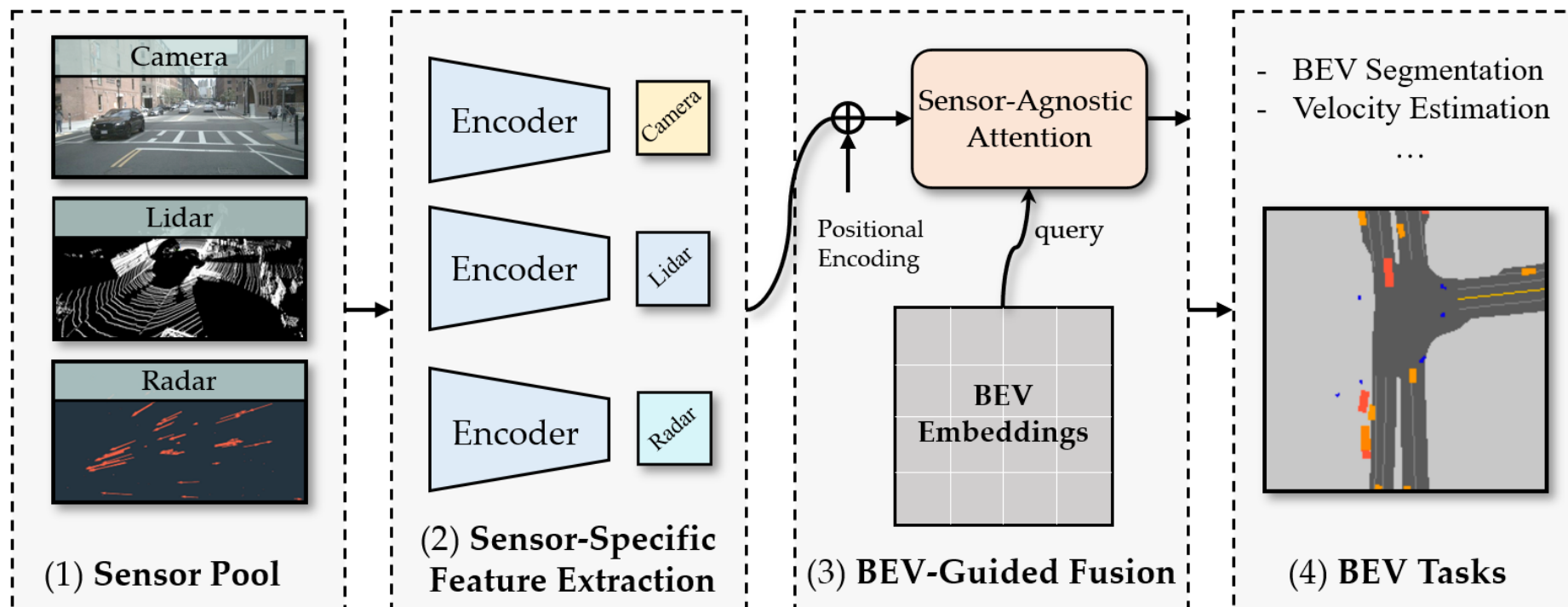
# Our Key Insight: Query Sensors from the Unified Space

Allow model to decide the **weights** of different sensors

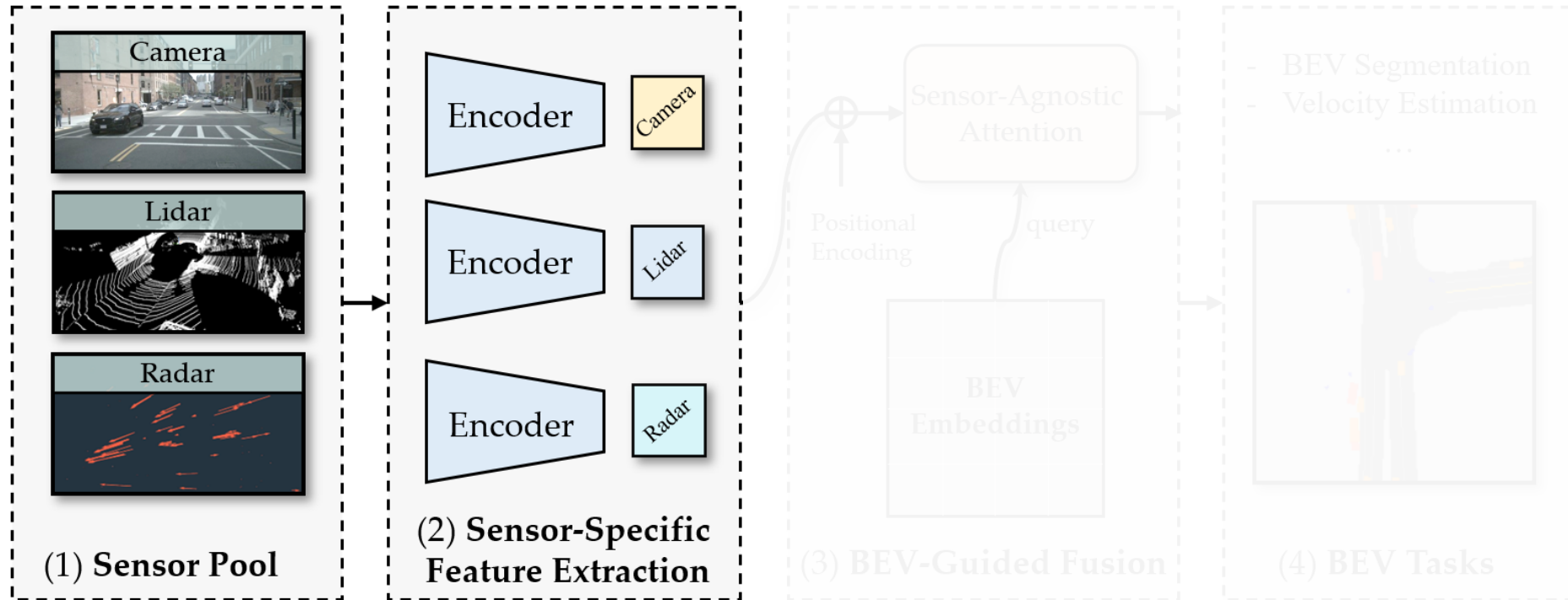
Build the representation **bottom-up** from BEV



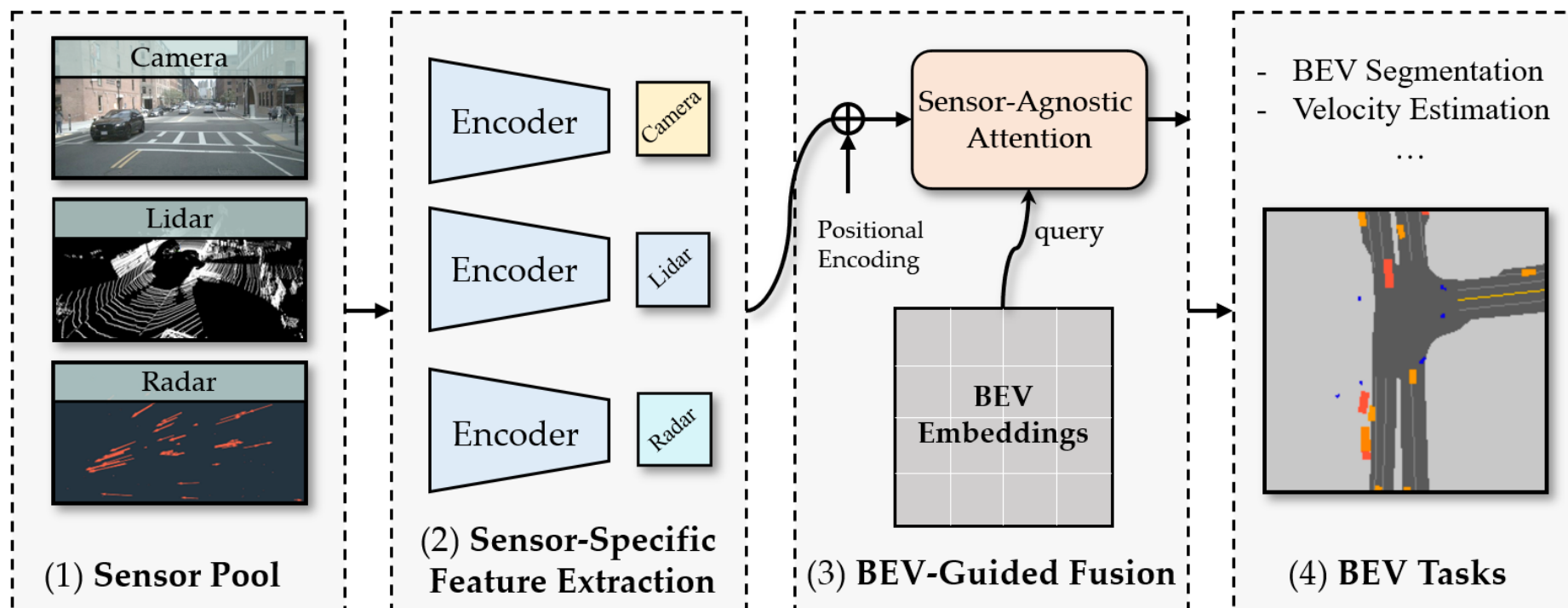
# Pipeline of BEVGuide



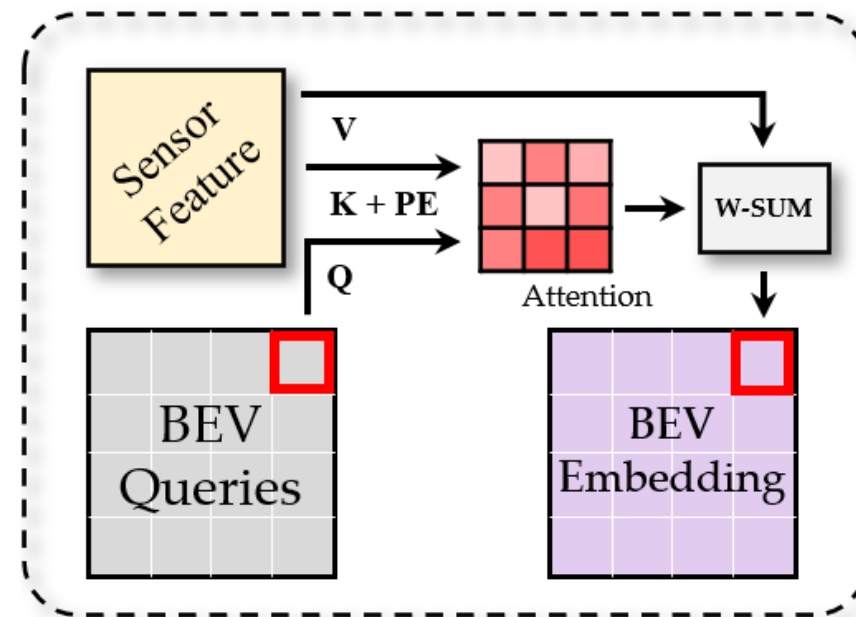
# Sensor Pool and Sensor-Specific Encoders



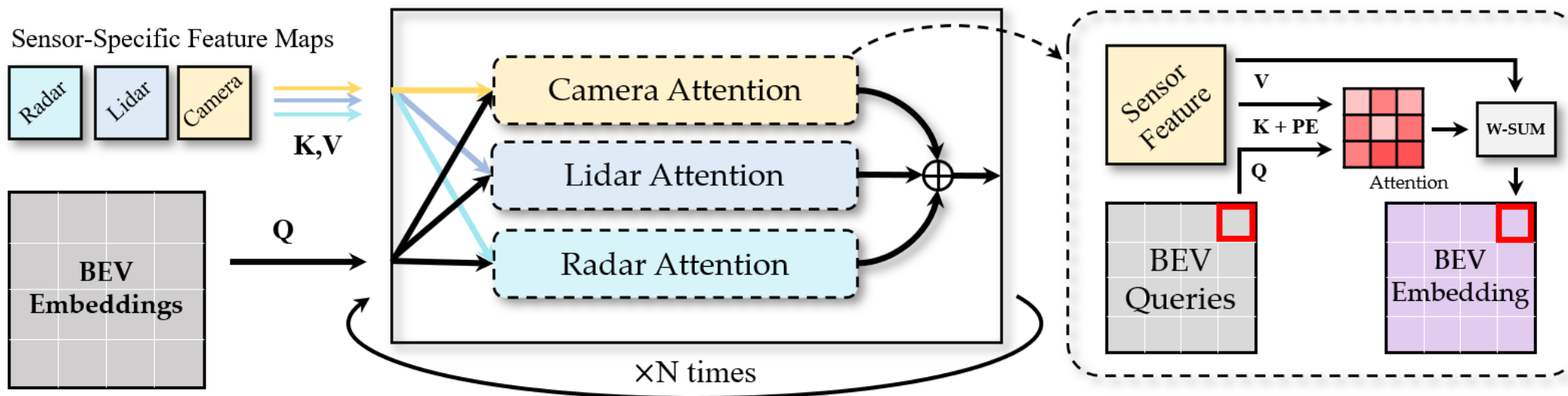
# BEV-Guided Sensor-Agnostic Attention



# BEV-Guided Sensor Agnostic Attention



# BEV-Guided Sensor Agnostic Attention



# Geometric Correspondence as Positional Embedding

- For Camera,

$$\begin{aligned}x^{(\text{im})} &\simeq KM X^{(\text{w})} \\ M^{-1}K^{-1}x^{(\text{im})} &\simeq X^{(\text{w})}\end{aligned}$$

Image Key                      BEV Query

$K$ : intrinsics

$M$ : extrinsics (pose)

# Geometric Correspondence as Positional Embedding

- For Camera,

$$\begin{aligned}x^{(\text{im})} &\simeq KMx^{(\text{w})} \\ M^{-1}K^{-1}x^{(\text{im})} &\simeq x^{(\text{w})}\end{aligned}$$

Image Key                      BEV Query

- For Lidar/Radar

$$\begin{aligned}x^{(\text{L/R})} &\simeq SX^{(\text{w})} \\ S^{-1}x^{(\text{L/R})} &\simeq X^{(\text{w})}\end{aligned}$$

Lidar/Radar Key                      BEV Query

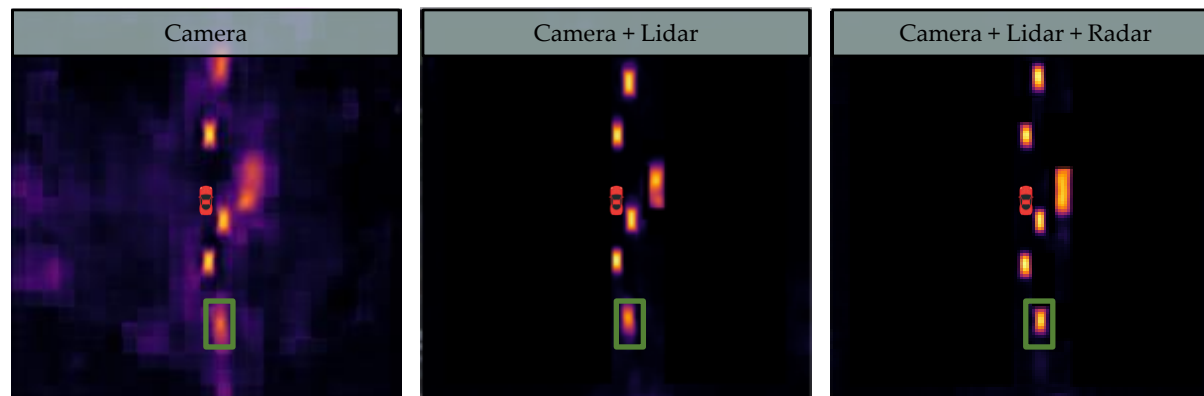
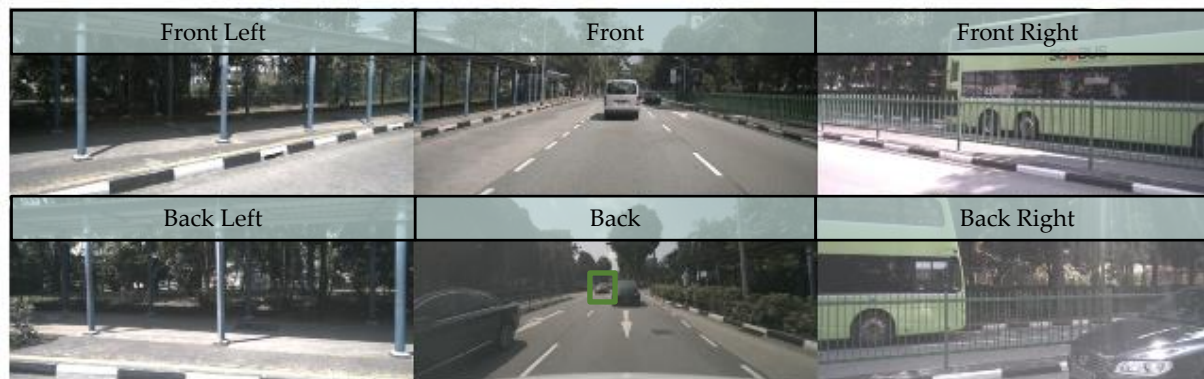
$K$ : intrinsics

$M$ : extrinsics (pose)

$S$ : scale & translate matrix



# Superior Performance on Segmentation



# Superior Performance on Diverse Driving Tasks: Segmentation, Detection, and Velocity Estimation

<i>Segmentation</i>	C R L	Vehicles	Roads
Cross-view	✓	36.0	74.3
FUTR3D	✓ ✓	46.6	-
Simple-BEV	✓ ✓	60.8	-
BEVFusion	✓ ✓	-	85.5
X-Align	✓ ✓	-	86.8
<b>BEVGuide</b>	✓ ✓ ✓	<b>79.0</b>	<b>86.9</b>

<i>Detection</i>	C R L	mAP	NDS
FUTR3D	✓ ✓	35.0	45.9
BEVGuide*	✓ ✓	42.1	53.7
BEVFusion	✓ ✓	68.5	71.4
BEVGuide*	✓ ✓	68.9	71.4
<b>BEVGuide</b>	✓ ✓ ✓	<b>69.3</b>	<b>71.5</b>

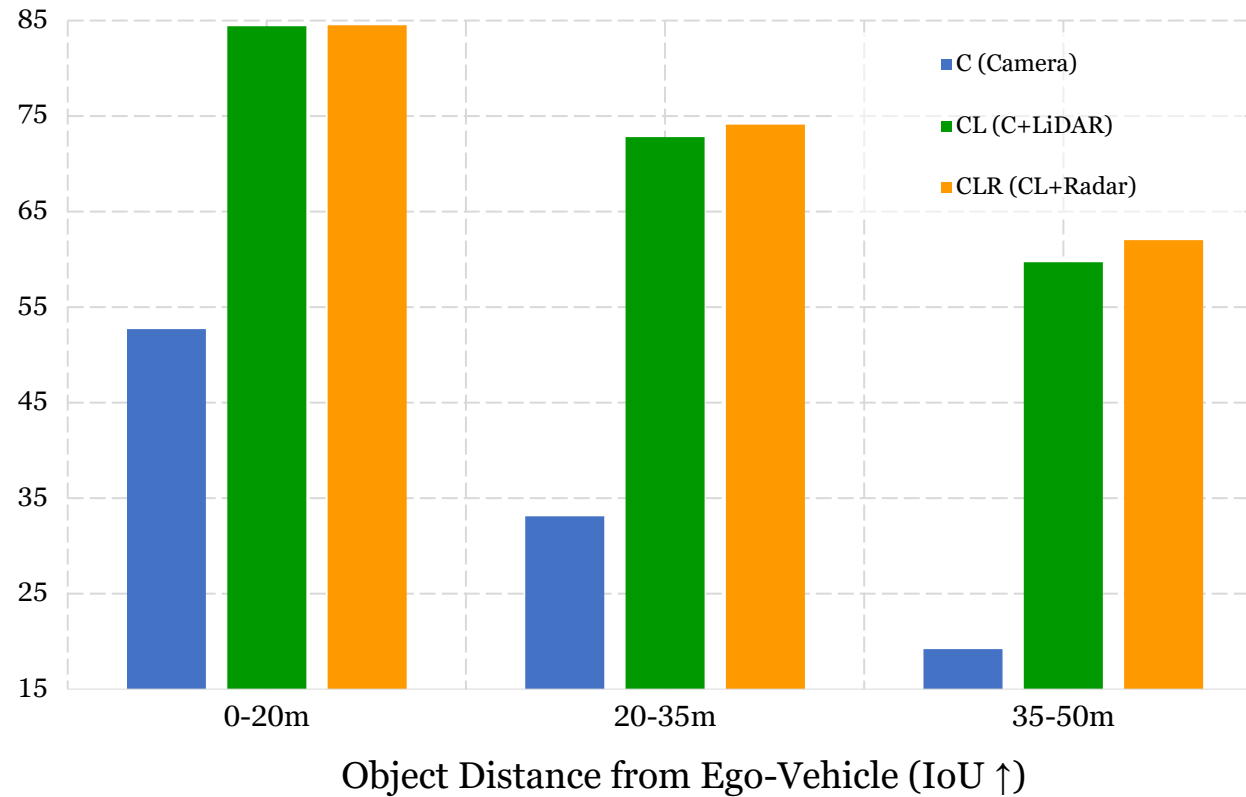
<i>Velo. Estimation</i>	C R L	P-AVE
Cross-view	✓	2.13
PointPainting	✓ ✓	1.90
BEVGuide*	✓ ✓	1.63
<b>BEVGuide</b>	✓ ✓ ✓	<b>0.81</b>

# Multi-Modality **Reduces** Domain Gap, and **Increases** Robustness

<i>Day → Night</i>	C R L	Day	Night	Gap
Cross-view	✓	40.4	18.8	21.6
BEVGuide*	✓ ✓	76.7	58.8	17.9
<b>BEVGuide</b>	<b>✓ ✓ ✓</b>	<b>79.5</b>	<b>64.2</b>	<b>15.3</b>

<i>Sunny → Rainy</i>	C R L	Day	Night	Gap
Cross-view	✓	37.3	28.1	9.2
BEVGuide*	✓ ✓	77.0	69.9	7.1
<b>BEVGuide</b>	<b>✓ ✓ ✓</b>	<b>80.7</b>	<b>74.6</b>	<b>6.1</b>

# Radar Sensor Helps Improve the Perception of More Distant Objects



# Summary

- BEVGuide, a **comprehensive** and **versatile** multi-modality fusion architecture
- Easily adapt to different **sensor combinations**
- Achieve **state-of-the-art** performance on various driving tasks

# Check Our Project Website



<https://yunzeman.github.io/BEVGuide/>