

Appendix

A Implementation Details

This section introduces the details utilized to implement LMD. First, we describe the dataset used for the BEV perception model employed to demonstrate our method, as well as the configuration of the fusion model utilized in the experiments.

A.1 Dataset

The nuScenes dataset [32] is a large-scale, multimodal dataset designed for autonomous driving tasks. It encompasses 1000 scenes, segmented into 700 for training, 150 for validation, and 150 for testing, each lasting about 20 seconds. The dataset captures a 360° horizontal field of view (FOV) through the use of six surround-view cameras, one lidar, and five radars, providing comprehensive environmental perception.

Notably, the nuScenes dataset stands out due to its high annotation frequency of every 0.5 seconds, translating to a 2 Hz rate, which includes over 1.4 million 3D bounding boxes across ten object categories. This extensive annotation facilitates advanced tasks like 3D object detection and tracking.

For evaluating detection performance, the dataset introduces unique metrics, including mean Average Precision (mAP) and five true positive (TP) metrics (ATE, ASE, AOE, AVE, and AAE) to measure errors in translation, scale, orientation, velocity, and attribute, respectively. These metrics are aggregated into the nuScenes Detection Score (NDS) to provide a comprehensive performance overview.

Furthermore, for the task like the Bird’s Eye View (BEV) segmentation, the dataset follows specific settings proposed in prior research, leveraging its rich radar point cloud data among other modalities. The nuScenes dataset not only facilitates the development of advanced autonomous driving technologies but also sets a benchmark for evaluating the performance of these systems through its detailed and nuanced evaluation metrics.

A.2 Model Architecture & Hyper-Parameters

The basic settings of the camera-radar fusion model follow those of SimpleBEV [22]. The camera branch processes input from six view images, each of which is passed through a backbone network, such as ResNet [36], to extract image features. These image features, along with radar point cloud information, are subsequently mapped onto the BEV coordinate system to facilitate perception tasks in the BEV space. The camera features and radar features in the BEV representation are concatenated and processed using convolutional operations to generate latent features. In this paper, this convolution process is defined as the fusion operation performed by the layer function ($l = 1$). The latent features are then passed through a U-Net [37] decoder, after which the prediction head carries out a segmentation task focused on identifying vehicle regions. Table 3 presents radar-camera and LiDAR-camera configuration details used for the SimpleBEV framework.

Table 3: Detailed configuration of modality-specific fusion setups

Parameter	Radar - Camera	LiDAR - Camera
Backbone	ResNet101	ResNet101
Fusion Operation	Concat & Conv	Concat & Conv
Sweeps	5	5
Use Meta Data	True	False
Input Size	(224, 400)	(224, 400)
BEV Coordinate	(200, 8, 200)	(200, 8, 200)

A.3 Experiments Compute Resources

We employed pre-trained encoders for camera, radar, and LiDAR modalities. Our method was evaluated for 6019 iterations with a batch size of 1 on 4 x NVIDIA RTX 3090 GPUs, which required approximately 3 hours in total.

B Proofs on LMD

B.1 Proposition 1: Handling Activation Layers

Let $A_j = F_j^{l-1}(x_c, x_r)$ denote the total input to the activation function for neuron j in layer l . Let $O_j = F_j^l(x_c, x_r)$ denote the original output of this activation function for neuron j in layer l . The linearized activation function for a modality-specific component h_{mj}^{l-1} is defined as $\hat{f}_j^l(h_{mj}^{l-1}) = c_j \cdot h_{mj}^{l-1}$. The constant c_j is given by $c_j = \frac{O_j}{A_j + \epsilon}$.

The total input to the activation function is the sum of its decomposed modality-specific components: $A_j = \sum_{m \in \{c, r, b\}} h_{mj}^{l-1}$. The output of the linearized layer $\hat{F}_j^l(x_c, x_r)$ is obtained by summing the outputs of the linearized activation function applied to each component (due to the linearity of $h \mapsto c_j h$ and as implied by Eq. (10) and (11) which state $h_{mj}^l = \hat{f}_j^l(h_{mj}^{l-1})$ and $\hat{f}_j^l(\sum_{m \in \{c, r, b\}} h_{mj}^{l-1}) = \sum_{m \in \{c, r, b\}} \hat{f}_j^l(h_{mj}^{l-1})$):

$$\hat{F}_j^l(x_c, x_r) = \sum_{m \in \{c, r, b\}} \hat{f}_j^l(h_{mj}^{l-1}) = \sum_{m \in \{c, r, b\}} c_j h_{mj}^{l-1} = c_j \sum_{m \in \{c, r, b\}} h_{mj}^{l-1} = c_j A_j.$$

Substituting the definition of c_j :

$$\hat{F}_j^l(x_c, x_r) = \left(\frac{O_j}{A_j + \epsilon} \right) A_j.$$

Eq (7) requires that the linearized model's output equals the original model's output at the operating point: $\hat{F}_j^l(x_c, x_r) = F_j^l(x_c, x_r)$, which is $\hat{F}_j^l(x_c, x_r) = O_j$. Therefore, we must have:

$$O_j = \left(\frac{O_j}{A_j + \epsilon} \right) A_j.$$

We analyze this equation in different cases:

- **Case 1:** $O_j \neq 0$. In this case, we can divide both sides by O_j :

$$1 = \frac{A_j}{A_j + \epsilon}.$$

This implies $A_j + \epsilon = A_j$, which means $\epsilon = 0$. Thus, if $O_j \neq 0$, eq. (7) is satisfied when $\epsilon = 0$.

- **Case 2:** $O_j = 0$. The equation becomes:

$$0 = \left(\frac{0}{A_j + \epsilon} \right) A_j.$$

This simplifies to $0 = 0$, provided that $A_j + \epsilon \neq 0$. This condition typically holds as ϵ is a small constant used to prevent division by zero. This case covers scenarios such as:

- $A_j = 0$ and $O_j = 0$ (e.g., ReLU(0)=0, GELU(0)=0). Here $c_j = 0/(0 + \epsilon) = 0$. Then $\hat{F}_j^l(x_c, x_r) = 0 \cdot 0 = 0$, which equals O_j . eq. (7) holds.
- $A_j \neq 0$ but $O_j = 0$ (e.g., ReLU applied to a negative input). Here $c_j = 0/(A_j + \epsilon) = 0$ (assuming $A_j + \epsilon \neq 0$). Then $\hat{F}_j^l(x_c, x_r) = 0 \cdot A_j = 0$, which equals O_j . eq. (7) holds.

The connection to DTD, where c_j is interpreted as the slope of the segment joining (A_j, O_j) (and possibly the origin, if the activation passes through it), underpins this choice. If $c_j = O_j/A_j$ (for $A_j \neq 0$), then $\hat{F}_j^l = (O_j/A_j)A_j = O_j$, directly satisfying eq. (7). The formulation with ϵ makes this robust.

Thus, the proposed construction for the linearized activation function \hat{f}_j^l ensures that eq. (7) is satisfied under the conditions discussed. \square

B.2 Proposition 2: Handling BatchNorm Layers

The linearized output of the layer is the sum over modality components:

$$\hat{F}_j^l(x_c, x_r) = \sum_{m \in \{c, r, b\}} \hat{f}_j^l(h_{mj}^{l-1}).$$

Camera ($m = c$) and radar ($m = r$) components. Because $\delta_c = \delta_r = 0$,

$$\hat{f}_j^l(h_{mj}^{l-1}) = \frac{\gamma_j^l}{\sqrt{(\sigma_j^l)^2 + \epsilon}} h_{mj}^{l-1} \quad (m \in \{c, r\}).$$

824 **Bias component** ($m = b$). With $\delta_b = 1$,

$$\hat{f}_j^l(h_{bj}^{l-1}) = \frac{\gamma_j^l}{\sqrt{(\sigma_j^l)^2 + \varepsilon}} (h_{bj}^{l-1} - \mu_j^l) + \beta_j^l.$$

825 Summing the three parts yields

$$\hat{F}_j^l(x_c, x_r) = \frac{\gamma_j^l}{\sqrt{(\sigma_j^l)^2 + \varepsilon}} (h_{cj}^{l-1} + h_{rj}^{l-1} + h_{bj}^{l-1} - \mu_j^l) + \beta_j^l.$$

826 Because the LMD framework guarantees $F_j^{l-1}(x_c, x_r) = \sum_{m \in \{c, r, b\}} h_{mj}^{l-1}$, we can rewrite the above as

$$\hat{F}_j^l(x_c, x_r) = \frac{\gamma_j^l}{\sqrt{(\sigma_j^l)^2 + \varepsilon}} (F_j^{l-1}(x_c, x_r) - \mu_j^l) + \beta_j^l,$$

827 which is the BatchNorm forward pass $F_j^l(x_c, x_r)$ in evaluation mode. Hence $\hat{F}_j^l = F_j^l$, satisfying eq. (7).

828 **B.3 Proposition 3: Handling InstanceNorm Layers**

829 Let $F^{l-1}(x_c, x_r)$ denote the input of the feature map to the InstanceNorm layer in layer l . For a specific channel
830 j within this feature map, let H_j^{l-1} represent the activations for the channel j . The standard InstanceNorm
831 operation for this channel is defined as:

$$F_j^l = \gamma_j \frac{H_j^{l-1} - \mathbb{E}[H_j^{l-1}]}{\sqrt{\text{Var}[H_j^{l-1}] + \epsilon}} + \beta_j$$

832 where $\mathbb{E}[H_j^{l-1}]$ and $\text{Var}[H_j^{l-1}]$ are the mean and variance of H_j^{l-1} computed over its spatial dimensions,
833 respectively. γ_j and β_j are learnable scaling and shifting parameters for channel j .

834 In the LMD framework, the input to the l -th layer, $F^{l-1}(x_c, x_r)$, is decomposed into modality-specific compo-
835 nents:

$$F^{l-1}(x_c, x_r) = \sum_{m \in \{c, r, b\}} h_m^{l-1}$$

836 where h_m^{l-1} represents the feature map attributed to the modality m (camera, radar, or bias) at layer $l - 1$. For
837 channel j , this means:

$$H_j^{l-1} = \sum_{m \in \{c, r, b\}} h_{mj}^{l-1}$$

838 The LMD framework requires that the linearized layer function \hat{f}^l satisfies the property that the sum of the
839 decomposed outputs equals the output of the linearized function applied to the sum of inputs, which in turn must
840 match the original function output at the operating point from eq. (7).

841 For InstanceNorm, we state that the variance $\sigma_j^l = \text{Var}[H_j^{l-1}]$ is treated as a constant, saved from a single
842 forward pass of the original fusion model.

843 By considering the proposed decomposition rule for $\hat{f}_j^l(h_{mj}^{l-1})$:

$$\hat{f}_j^l(h_{mj}^{l-1}) = \frac{h_{mj}^{l-1} - \mathbb{E}[h_{mj}^{l-1}]}{\sqrt{\sigma_j^l + \epsilon}} \gamma_j^l + \delta_m \beta_j^l$$

844 To sum these decomposed outputs over all modalities $m \in \{c, r, b\}$ to see if they reconstruct the original
845 InstanceNorm output F_j^l :

$$\sum_{m \in \{c, r, b\}} \hat{f}_j^l(h_{mj}^{l-1}) = \sum_{m \in \{c, r, b\}} \left(\frac{h_{mj}^{l-1} - \mathbb{E}[h_{mj}^{l-1}]}{\sqrt{\sigma_j^l + \epsilon}} \gamma_j^l + \delta_m \beta_j^l \right)$$

846 By splitting the summation:

$$= \frac{\gamma_j^l}{\sqrt{\sigma_j^l + \epsilon}} \sum_{m \in \{c, r, b\}} (h_{mj}^{l-1} - \mathbb{E}[h_{mj}^{l-1}]) + \sum_{m \in \{c, r, b\}} \delta_m \beta_j^l$$

For the first term, we use the linearity of expectation:

$$\begin{aligned} \sum_{m \in \{c, r, b\}} (h_{mj}^{l-1} - \mathbb{E}[h_{mj}^{l-1}]) &= \sum_{m \in \{c, r, b\}} h_{mj}^{l-1} - \sum_{m \in \{c, r, b\}} \mathbb{E}[h_{mj}^{l-1}] \\ &= H_j^{l-1} - \mathbb{E} \left[\sum_{m \in \{c, r, b\}} h_{mj}^{l-1} \right] \\ &= H_j^{l-1} - \mathbb{E}[H_j^{l-1}] \end{aligned}$$

For the second term, within the definition of δ_m : $\delta_m = 0$ for $m \in \{c, r\}$ and $\delta_m = 1$ for $m \in \{b\}$. Therefore:

$$\sum_{m \in \{c, r, b\}} \delta_m \beta_j^l = \delta_c \beta_j^l + \delta_r \beta_j^l + \delta_b \beta_j^l = 0 \cdot \beta_j^l + 0 \cdot \beta_j^l + 1 \cdot \beta_j^l = \beta_j^l$$

847 Substituting these back into the summed expression:

$$\sum_{m \in \{c, r, b\}} \hat{f}_j^l(h_{mj}^{l-1}) = \frac{\gamma_j^l (H_j^{l-1} - \mathbb{E}[H_j^{l-1}])}{\sqrt{\sigma_j^l + \epsilon}} + \beta_j^l$$

848 Since σ_j^l is taken as the fixed variance $\text{Var}[H_j^{l-1}]$ from the first forward pass, this expression is the original
 849 InstanceNorm output F_j^l . Since each modality is processed independently except for the shared, fixed scale, the
 850 resulting features adhere to the intended *separation property*. The proposed decomposition rule for InstanceNorm
 851 layers ensures that the sum of the modality-specific outputs from the linearized layer perfectly reconstructs the
 852 output of the original InstanceNorm layer, satisfying eq. (7) of the LMD framework.

853 C Variants on Activation Layers

854 Necessities of considering the bias-splitting rules in activation layers, aside from cases where a constant term is
 855 inevitably introduced (e.g., normalization layers), is as follows. When a model is linearized, a bias' prediction
 856 can be obtained by inputting zeros into all data points. As shown in fig. 2, it is evident that bias features also
 857 exhibit perception capabilities. Therefore, the process of appropriately distributing the bias contribution between
 858 the camera and radar can be considered. Particularly in the activation layers, if a specific neuron is activated by
 859 the camera, adding a positive bias to the camera feature at that neuron could enhance the camera's contribution.
 860 Conversely, if a neuron is activated by the radar, adding a positive bias to the radar feature could strengthen the
 861 radar's contribution. This idea led to the development of the sum rule introduced in appendix C.1. The modality
 862 predictions based on various bias-splitting rules need to be extensively explored in future experiments.

863 While the sum rule adds the entire bias to whichever modality actually triggers a ReLU, it disregards how strongly
 864 the two modalities contribute in magnitude. Empirically this may still blur modality separation when both camera
 865 and radar deliver noticeable—but differently scaled—signals. The ratio splitting rule therefore apportions the
 866 bias term h_{bj}^{l-1} proportionally to the absolute pre-activations of camera and radar. The impact of sum and ratio
 867 splitting rules can be assessed via the four perturbation-based metrics (R_p/R , R_p/C , C_p/R , C_p/C). These
 868 rules modulate how sensitively and independently each modality responds to perturbation, offering an empirical
 869 basis for bias handling in activation layers.

870 C.1 Sum Splitting Rule

871 For the most common activation function, ReLU, if a camera feature value is less than zero but a radar feature
 872 value is greater than zero, resulting in the activation of a specific neuron, it is reasonable to interpret that the
 873 neuron was activated by the radar. In such cases, a splitting rule can be applied where the bias feature is added to
 874 the radar feature when the bias feature value is greater than zero. Conversely, if both a camera feature value
 875 and a bias feature value are greater than zero while the radar feature value is less than zero, distributing the bias
 876 feature to the camera feature strengthen the camera's contribution to that neuron. The sum condition tested in
 877 our experiments are derived from this idea. The related equations are as follows:

$$\begin{aligned} \text{Cam} - \text{Condition} &= h_{cj}^{l-1} < 0, h_{rj}^{l-1} > 0, h_{bj}^{l-1} < 0 \\ &\text{and } h_{cj}^{l-1} > 0, h_{rj}^{l-1} < 0, h_{bj}^{l-1} > 0 \end{aligned}$$

$$\begin{aligned} \text{Rad} - \text{Condition} &= h_{cj}^{l-1} < 0, h_{rj}^{l-1} > 0, h_{bj}^{l-1} > 0 \\ &\text{and } h_{cj}^{l-1} > 0, h_{rj}^{l-1} < 0, h_{bj}^{l-1} < 0 \end{aligned}$$

$$\begin{aligned} \hat{f}^l(h_{cj}^{l-1}) &= \begin{cases} c(h_{cj}^{l-1} + h_{bj}^{l-1}), & \text{for Cam} - \text{Condition}, \\ c(h_{cj}^{l-1}), & \text{otherwise.} \end{cases} \\ \hat{f}^l(h_{rj}^{l-1}) &= \begin{cases} c(h_{rj}^{l-1} + h_{bj}^{l-1}), & \text{for Rad} - \text{Condition}, \\ c(h_{rj}^{l-1}), & \text{otherwise.} \end{cases} \\ \hat{f}^l(h_{bj}^{l-1}) &= \begin{cases} 0, & \text{for Cam \& Rad} - \text{Condition}, \\ c(h_{bj}^{l-1}), & \text{otherwise.} \end{cases} \\ , \text{ where } c &= \frac{F_j^l(\mathbf{x}_c, \mathbf{x}_r)}{F_j^{l-1}(\mathbf{x}_c, \mathbf{x}_r) + \epsilon}. \end{aligned} \tag{13}$$

880 C.2 Ratio Splitting Rule

881 In activation layers a bias term h_{bj}^{l-1} can enhance the modality responsible for the neuron's activation. The ratio
882 splitting rule redistributes that bias to the camera and radar features in proportion to their absolute magnitudes,
883 while leaving neurons that fail the ratio conditions unchanged. This preserves the conservation property in eq. (7)
884 and generalizes the identity and uniform rules introduced earlier.

$$\begin{aligned} \text{Ratio} - \text{Condition} &= (h_{cj}^{l-1} > 0, h_{rj}^{l-1} > 0, h_{bj}^{l-1} > 0) \\ &\text{or } (h_{cj}^{l-1} < 0, h_{rj}^{l-1} < 0, h_{bj}^{l-1} < 0), \\ \text{Ratio} - \text{Condition2} &= (h_{cj}^{l-1} > 0, h_{rj}^{l-1} > 0, h_{bj}^{l-1} < 0) \\ &\text{or } (h_{cj}^{l-1} < 0, h_{rj}^{l-1} < 0, h_{bj}^{l-1} > 0). \end{aligned}$$

$$\alpha_j = \frac{|h_{rj}^{l-1}|}{|h_{cj}^{l-1}| + |h_{rj}^{l-1}| + \varepsilon}, \quad \alpha_j \in [0, 1],$$

885 with a small $\varepsilon > 0$ for numerical stability.

886 Let $c = \frac{F_j^l(x_c, x_r)}{F_j^{l-1}(x_c, x_r) + \varepsilon}$ be the slope from proposition 1. Then

$$\begin{aligned} \hat{f}^l(h_{cj}^{l-1}) &= \begin{cases} \alpha(h_{cj}^{l-1} + (1 - \alpha_j)h_{bj}^{l-1}), & \text{Ratio} - \text{Condition}, \\ \alpha(h_{cj}^{l-1} + \alpha_j h_{bj}^{l-1}), & \text{Ratio} - \text{Condition2}, \\ c h_{cj}^{l-1}, & \text{otherwise,} \end{cases} \\ \hat{f}^l(h_{rj}^{l-1}) &= \begin{cases} \alpha(h_{rj}^{l-1} + \alpha_j h_{bj}^{l-1}), & \text{Ratio} - \text{Condition}, \\ \alpha(h_{rj}^{l-1} + (1 - \alpha_j)h_{bj}^{l-1}), & \text{Ratio} - \text{Condition2}, \\ c h_{rj}^{l-1}, & \text{otherwise,} \end{cases} \\ \hat{f}^l(h_{bj}^{l-1}) &= \begin{cases} 0, & \text{Ratio} - \text{Condition or Ratio} - \text{Condition2}, \\ c h_{bj}^{l-1}, & \text{otherwise.} \end{cases} \end{aligned}$$

887 Because h_{bj}^{l-1} is only re-partitioned, the eq. (7) is strictly preserved.

888 The ratio rule thus interpolates smoothly between existing bias-handling strategies while adapting to the
889 magnitude of each modality's evidence, leading to the stability improvements reported in section 5.

Table 4: Comparison of Pearson Correlation and Mean Squared Error for selected LMD rules (Radar + Camera)

Modality	Method	Pearson Correlation				Mean Squared Error			
		R _p /R (↓)	R _p /C (↑)	C _p /R (↑)	C _p /C (↓)	R _p /R (↑)	R _p /C (↓)	C _p /R (↓)	C _p /C (↑)
Radar + Camera	Sum-Splitting	0.12 ± 0.04	1.00 ± 0.00	1.00 ± 0.00	0.23 ± 0.05	9.63 ± 2.34	0.00 ± 0.00	0.00 ± 0.00	40.43 ± 19.35
	Ratio-Splitting	0.08 ± 0.06	1.00 ± 0.00	1.00 ± 0.00	0.18 ± 0.04	4.46 ± 1.45	0.00 ± 0.00	0.00 ± 0.00	30.35 ± 4.19

D Further Experiments

D.1 Further Post-hoc Interpretation

This section presents additional examples of applying LMD to the BEV perception task. Through further visualizations, we expect that our proposed method is validated with sufficient consistency to reliably interpret the model.

fig. 3 and fig. 4 show further visualization results for the reliability of the proposed LMD on the BEV perception benchmark. fig. 3 normalizes each modality-specific feature map with a sigmoid transform before projection, so only positive activations remain; the bright regions therefore mark locations where a single sensor delivers decisive, supportive cues to the detector, making the saliency of true objects immediately apparent. By contrast, fig. 4 retains the full signed output and colour-codes positive and negative responses, allowing us to see not only where a modality reinforces the final prediction but also where it actively suppresses prediction.

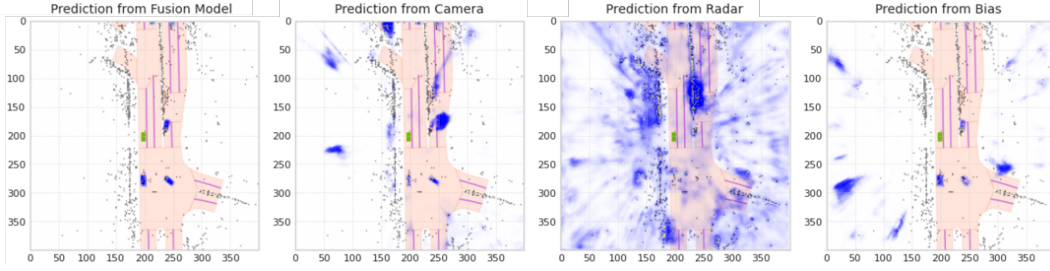


Figure 3: Visualizations of using Sigmoid.

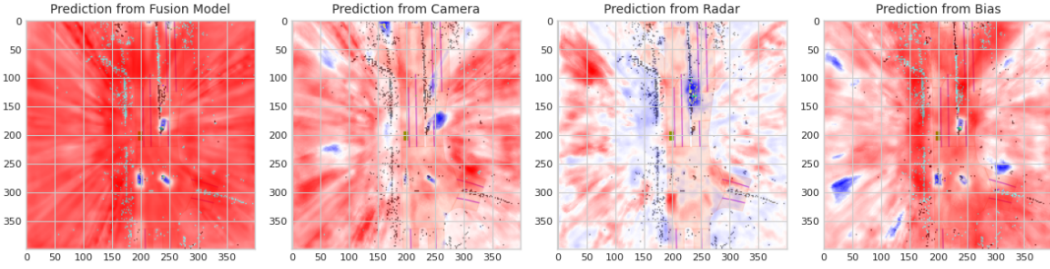


Figure 4: Visualizations of both Positive and Negative Values.

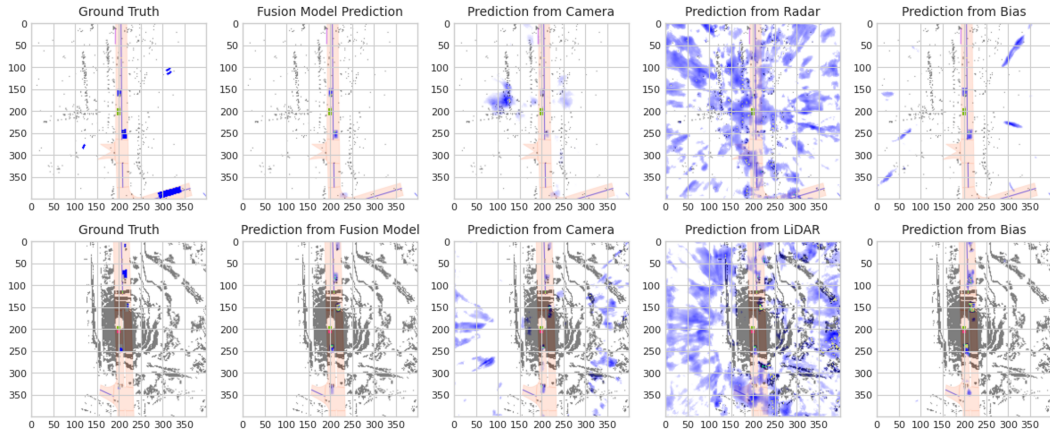


Figure 5: Further Qualitative Visualization Results

D.2 Metric Formulas

To evaluate the separation property of the modality-specific decomposition, we introduce four perturbation-based metrics. Here, perturbation refers to the replacement of one modality input with an uncorrelated sample. These metrics assess how much each modality-specific output responds to changes in its own input and remains invariant to others. Each metric is computed using the Pearson Correlation Coefficient (PCC) and Mean Squared Error (MSE).

Pearson Correlation Coefficient (PCC). Briefly, for pearson correlation coefficient, we define:

$$\text{PCC}_{k,m} = \frac{\sum_i (P_k[i] - \bar{P}_k) (P_{(k+500m) \bmod N}[i] - \bar{P}_{(k+500m) \bmod N})}{\sqrt{\sum_i (P_k[i] - \bar{P}_k)^2 \sum_i (P_{(k+500m) \bmod N}[i] - \bar{P}_{(k+500m) \bmod N})^2}}, \quad (14)$$

where $N = 6019$, P_k and $P_{(k+500m) \bmod N}$ are the predictions for the sample spaced by $500m$ from k , \bar{P}_k and $\bar{P}_{(k+500m) \bmod N}$ are their respective mean values, and $m \in \{1, 2, \dots, 12\}$ ensures 12 distinct comparisons. PCC measures the similarity between two prediction vectors by first subtracting their respective means (mean-centering), computing the dot product of the centered vectors, and normalizing it by the product of their standard deviations. The resulting value lies in the range $[-1, 1]$, where 1 indicates perfect correlation, 0 indicates no correlation, and -1 indicates perfect inverse correlation.

Mean Squared Error (MSE). In addition to PCC, we compute the mean squared error to quantify absolute differences:

$$\text{MSE}_{k,m} = \frac{1}{M} \sum_{i=1}^M (P_k[i] - P_{(k+500m) \bmod N}[i])^2, \quad (15)$$

where M denotes the total number of pixels in the BEV feature map.

Perturbation-based Metrics. Let x_c and x_r denote the clean camera and radar inputs, and x_c^{pert} , x_r^{pert} their respective perturbed versions. Let $F_{\text{camera}}(x_c, x_r)$ and $F_{\text{radar}}(x_c, x_r)$ denote the modality-specific predictions.

We define the following four metrics:

- R_p/R . Radar prediction change under radar perturbation:

$$\text{PCC}_{k,m}(F_{\text{radar}}(x_c, x_r^{\text{pert}}), F_{\text{radar}}(x_c, x_r)), \quad \text{MSE}_{k,m}(F_{\text{radar}}(x_c, x_r^{\text{pert}}), F_{\text{radar}}(x_c, x_r))$$

Lower PCC and higher MSE are better, indicating sensitivity to radar input.

- R_p/C . Camera prediction change under radar perturbation:

$$\text{PCC}_{k,m}(F_{\text{camera}}(x_c, x_r^{\text{pert}}), F_{\text{camera}}(x_c, x_r)), \quad \text{MSE}_{k,m}(F_{\text{camera}}(x_c, x_r^{\text{pert}}), F_{\text{camera}}(x_c, x_r))$$

Higher PCC and lower MSE are better, indicating invariance of the camera modality.

- C_p/R . Radar prediction change under camera perturbation:

$$\text{PCC}_{k,m}(F_{\text{radar}}(x_c^{\text{pert}}, x_r), F_{\text{radar}}(x_c, x_r)), \quad \text{MSE}_{k,m}(F_{\text{radar}}(x_c^{\text{pert}}, x_r), F_{\text{radar}}(x_c, x_r))$$

Higher PCC and lower MSE are better, indicating invariance of the radar modality.

- C_p/C . Camera prediction change under camera perturbation:

$$\text{PCC}_{k,m}(F_{\text{camera}}(x_c^{\text{pert}}, x_r), F_{\text{camera}}(x_c, x_r)), \quad \text{MSE}_{k,m}(F_{\text{camera}}(x_c^{\text{pert}}, x_r), F_{\text{camera}}(x_c, x_r))$$

Lower PCC and higher MSE are better, indicating sensitivity to camera input.

These metrics are averaged over the entire test set using uniformly sampled perturbations. Together, they verify whether the decomposition satisfies: (1) sensitivity to the perturbed modality and (2) invariance to the unperturbed modality — key indicators of successful modality separation.

E Discussion

E.1 Broader Impact

LMD introduces the first post-hoc framework capable of attributing the predictions of each layer in a multi-sensor fusion model to its respective modalities, without modifying the original architecture.

Positive societal impacts may include:

936 • **Enhanced safety in autonomous systems.** By revealing the sensor modality (camera, LiDAR, or
937 radar) primarily responsible for each prediction, LMD enables identification of single-sensor failures
938 before they propagate into critical errors.

939 • **Accelerated certification and debugging.** Modality-level attributions can support system audits and
940 facilitate the generation of safety evidence, in line with emerging regulatory requirements for ADAS
941 systems.

942 • **Extension of interpretability to other multi-modal domains.** Since LMD is agnostic to model
943 architecture, it can be applied to a wide range of multimodal perception systems beyond autonomous
944 driving, such as medical image–text fusion or surveillance video analysis, to clearly attribute the role
945 of each input modality in the model’s decision-making process.

946 • **Commitment to open science.** We will publicly release our codebase, pretrained checkpoints, and
947 evaluation scripts under an open-source license to promote reproducibility and rigorous external
948 validation.

949 **Potential risks include:**

950 • **Misleading reassurance.** Users may over-trust visual explanations; a clear modality-wise attribution
951 does not guarantee that the fused decision is accurate.

952 • **Bias amplification.** LMD uncovers but does not correct dataset imbalances, and naive use of its
953 explanations may perpetuate structural biases.

954 • **Residual opacity.** The interpretability of high-order interactions and root-point selections in activation
955 layers remains an open challenge for future work.

956 To mitigate these risks, we recommend pairing LMD with the perturbation-based metrics introduced in this paper
957 to quantitatively assess explanation fidelity prior to deployment. We also encourage conducting independent
958 reproducibility audits on geographically and demographically diverse datasets, and suggest applying the same
959 level of scrutiny—augmented by domain-specific expertise and operational context—when extending LMD
960 to other multimodal domains, such as medical imaging or service robotics. LMD sets a new benchmark for
961 interpretability in multimodal domains and provides a solid foundation for future research and real-world
962 deployment.

963 **E.2 Model-Wise Explanation**

964 Understanding the extent to which we can perform model-wise explanations by comparing the camera-only
965 model and the fusion model is a critical issue. Interpreting the contributions of each modality based on the
966 results of these two models for specific data is highly challenging. This difficulty arises, for the two models are
967 trained on different input distributions.

968 For instance, if a fusion model successfully detects an object that camera-only model fails to detect, it would be
969 a clear error to attribute this success to radar data. Predictions from the camera-only model should be utilized
970 only for references. For example, if the camera-only model successfully detects a particular vehicle, it can be
971 inferred that the information captured by the camera sensor for that region is sufficient. This insight can be
972 considered when analyzing the fusion model using the LMD approach.

973 Nonetheless, it is crucial to recognize that LMD is designed to provide post-hoc interpretations for a single
974 trained model, not to be a tool for interpreting two different models trained on different data distributions.

975 **E.3 GradCAM Experiments**

976 Since LMD enables access to internal activation maps, it is possible to compute activation maps from radar data
977 and subsequently apply the Grad-CAM [38] approach.

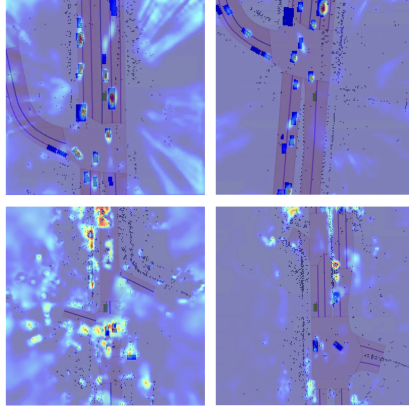


Figure 6: Visualizations of Radar-GradCAM Using Radar's Activations from LMD.

978 Grad-CAM [38] leverages gradient information to highlight relevant areas, Grad-CAM++ [39] introduces
 979 modified approach by taking positive gradients of target functions, and Score-CAM [40] bypasses the need for
 980 gradient information altogether, instead relying on activation scores to determine feature significance. Seg-Grad-
 981 CAM [41] is an adaptation of the Grad-CAM method, applied to the task of image segmentation by treating the
 982 segmentation output as the target function while maintaining the original mechanism of Grad-CAM for visual
 983 explanations. This approach enables a focused analysis of regions relevant to segmentation tasks, leveraging the
 984 gradients of the target output with respect to the feature maps. The Grad-CAM is calculated by :

$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\sum_k \alpha_k^c A^{lk} \right) \quad (16)$$

985 where $L_{\text{Grad-CAM}}^c$ is the class activation map for class c , α_k^c are the weights for the k -th feature map in l -th layer
 986 A^{lk} , computed as the global average of the gradients flowing back from the output unit for class c . ReLU is
 987 applied to focus on features that have a positive influence on the class of interest. The weight of k -th feature
 988 map for class c is computed as follows.

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^{lk}} \quad (17)$$

989 where y^c , A_{ij}^{lk} , and Z denote the score for class c , (i, j) -th entry in k -th activation map, and normalization factor
 990 respectively. In Seg-Grad-CAM, y^c is replaced by $\sum_{(i,j) \in R} y_{ij}^c$, where R is a set of pixel indices of interest in
 991 the output mask.

992 Through obtained activations which are derived solely from the radar data with LMD, it becomes possible to use
 993 the conventional CAM method based on the activations at the intermediate layer, which we call Radar-GradCAM.
 994 The Radar-GradCAM is calculated by :

$$L_{\text{Radar-GradCAM}}^c = \text{ReLU} \left(\sum_k \alpha_k^c (\hat{A}^{lk}(\mathbf{o}_c, \mathbf{x}_r) - \hat{A}^{lk}(\mathbf{o}_c, \mathbf{o}_r)) \right) \quad (18)$$

995 , where $\hat{A}^{lk}(\mathbf{o}_c, \mathbf{x}_r)$ and $\hat{A}^{lk}(\mathbf{o}_c, \mathbf{o}_r)$ are the activations of linearized fusion model \hat{F} and α_k^c is calculated in
 996 the same manner as in eq. (17).

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial \sum_{(i,j) \in R} \hat{F}_j^l(\mathbf{x}_c, \mathbf{x}_r)}{\partial A_{ij}^{lk}(\mathbf{x}_c, \mathbf{x}_r)} \quad (19)$$

997 Based on this calculation, the Radar-GradCAM shown in fig. 6 reveals which parts of the fusion model are
 998 relevant to the radar data points, by taking the radar-only activations that are crucial to the perception task.

999 F Extention to General Multimodal Model

1000 In this Section, we demonstrate that the LMD method can be generalized to multimodal models. We achieve this
1001 by extending the formulation to accommodate M modalities

1002 Beginning at the fusion layer ($l = 1$), the first-order Taylor expansion of f_j^1 becomes :

$$\begin{aligned} f_j^1(\mathbf{x}_{m_1}, \dots, \mathbf{x}_{m_M}) &= \sum_i \sum_{m \in \{m_1, \dots, m_M\}} J_{mji}^1 x_{mi} + \\ &\underbrace{\left(f_j^1(\tilde{x}_{m_1}, \dots, \tilde{x}_{m_M}) - \sum_i \sum_{m \in \{m_1, \dots, m_M\}} J_{mji} \tilde{x}_{mi} + \epsilon \right)}_{b_j^1} \end{aligned} \quad (20)$$

1003 where m indexes the modality, $\tilde{x}_{m_1}, \dots, \tilde{x}_{m_M}$ represent reference points for M modalities, and \mathbf{J}_{mji} denotes
1004 the Jacobian calculated at the reference points for modality m .

1005 For subsequent layer with $l = 2, \dots, N$, suppose l -th layer function takes $\sum_{m \in \{m_1, \dots, m_M, b\}} h_{mj}^{l-1}$ as inputs,
1006 the following holds :

$$\begin{aligned} f_j^l \left(\sum_{m \in \{m_1, \dots, m_M, b\}} h_{mj}^{l-1} \right) &= \\ \sum_i \sum_{m \in \{m_1, \dots, m_M, b\}} J_{mji}^{l-1} (h_{mi}^{l-1}) + b_j^l &= \sum_{m \in \{m_1, \dots, m_M, b\}} h_{mj}^l \end{aligned} \quad (21)$$

1007 We define all modality features that all layers in fusion model be decomposed into them.

$$F_j^l(\mathbf{x}_{m_1}, \dots, \mathbf{x}_{m_M}) = f_j^l \left(\sum_{m \in \{m_1, \dots, m_M, b\}} h_{mj}^{l-1} \right) = h_{m_1j}^l + \dots + h_{m_Mj}^l + h_{bj}^l \quad (22)$$

1008 From eq. (22), we can confirm inductively that the following holds. For the linearized layer functions $\hat{f}^1, \dots, \hat{f}^N$
1009 with bias-splitting rules applied, $F_j^l(\mathbf{x}_{m_1}, \dots, \mathbf{x}_{m_M}) = \hat{F}_j^l(\mathbf{x}_{m_1}, \dots, \mathbf{x}_{m_M})$ for $l \in \{1, \dots, N\}$, where
1010 $\hat{F}^l(\mathbf{x}_{m_1}, \dots, \mathbf{x}_{m_M})$ be $\hat{f}^l \circ \dots \circ \hat{f}^2 \circ \hat{f}^1(\mathbf{x}_{m_1}, \dots, \mathbf{x}_{m_M})$, forcing the epsilon in eq. (20) to zero. Based on
1011 all the layer functions, $\hat{f}^1, \dots, \hat{f}^N$, the modality features are computed as follows :

$$h_{mj}^l = \hat{f}_j^l(h_{mj}^{l-1}), m \in \{m_1, \dots, m_M, b\} \quad (23)$$

1012 Furthermore, the following arithmetic is satisfied for layer function \hat{f}_j^l :

$$\begin{aligned} \hat{f}_j^l \left(\sum_{m \in \{m_1, \dots, m_M, b\}} h_{mj}^{l-1} \right) &= \\ = \hat{f}_j^l(h_{m_1j}^{l-1}) + \dots + \hat{f}_j^l(h_{m_Mj}^{l-1}) + \hat{f}_j^l(h_{bj}^{l-1}) \end{aligned} \quad (24)$$

1013 From eq. (23) and eq. (24), the following property holds for every layers :

$$\begin{aligned} \hat{F}_j^l(\mathbf{x}_{m_1}, \dots, \mathbf{x}_{m_M}) &= \underbrace{\hat{F}_j^l(\mathbf{x}_{m_1}, \mathbf{o}_{m_2}, \dots, \mathbf{o}_{m_M})}_{h_{m_1j}^l} \\ &+ \dots + \underbrace{\hat{F}_j^l(\mathbf{o}_{m_1}, \dots, \mathbf{o}_{m_{M-1}}, \mathbf{x}_{m_M})}_{h_{m_Mj}^l} \\ &+ \underbrace{\hat{F}_{m_M}^l(\mathbf{o}_{m_1}, \dots, \mathbf{o}_{m_M})}_{h_{bj}^l}. \end{aligned} \quad (25)$$

1014 The linearization and bias-splitting rules applied to the normalization and activation layers are consistent with
1015 those described in the main paper.