# Foundation Models and GenAI (CSL7860)

## MidSem Exam
### *SED: A Simple Encoder-Decoder for Open-Vocabulary Semantic Segmentation*

Soumen Kumar
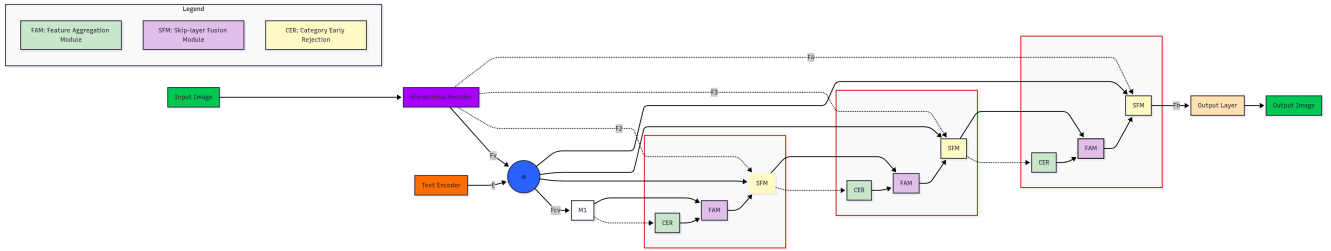Roll No-B22ES006

# 1 Executive Snapshot of the Paper

The SED paper aims to bridge the efficiency–accuracy gap in open-vocabulary semantic segmentation by combining a hierarchical convolutional image encoder with a simple, gradual fusion decoder. The result is a model that achieves strong open-vocabulary segmentation performance at a fraction of the computational cost of transformer-based or two-stage systems. Below, each core aspect is summarized.

## 1.1 Problem & Motivation

Most approaches to open-vocabulary segmentation are either too costly (transformers with quadratic complexity) or inefficient due to separate modules for mask and classification. SED addresses the urgent need for models that can deliver open-vocabulary segmentation with both speed and accuracy as large-scale vision-language datasets become mainstream.

## 1.2 Core Idea & Design

SED replaces the transformer backbone with a ConvNeXt encoder, producing multi-scale features for efficient and expressive pixel-to-class matching. A simple cosine similarity generates a dense cost map between image features and frozen CLIP text embeddings. The decoder then performs gradual, stage-wise fusion of cost maps with features using large-kernel convolutions, skip connections, and category early rejection to accelerate inference without loss.



## 1.3 Main Contributions

- An encoder-decoder framework that uses a hierarchical backbone for better efficiency and local context.

- A pixel-wise cost map based on cosine similarity between dense visual and text features, supporting open-vocabulary segmentation.

- A gradual fusion decoder design, with category early rejection, enabling up to 4.7× faster inference.

- SOTA performance on multiple datasets with faster average inference—e.g., 31.6% mIoU on ADE20K A-150 at 82 ms latency.

- Ablation studies demonstrating observable gains from each component, especially hierarchical encoding and gradual fusion.

## 1.4 Takeaway and Limitation

In essence, SED demonstrates the power of simplicity: a well-chosen hierarchical encoder and tactical gradual fusion decode suffice for robust, efficient open-vocabulary segmentation. The most salient limitation is SED's occasional confusion between near-synonym classes (e.g., "sea" versus "water"). Further research on fine-grained visual-language distinctions or attention strategies may lessen this issue in future work.

# 2 Method Deep-Dive

## 2.1 Mechanics: Architecture, Training Objective, Decoding Pipeline

The system proposed in the paper comprises three interconnected components working synergistically to achieve efficient pixel-level classification across arbitrary semantic categories.

### Hierarchical Encoder-based Cost Map Generation (HECG)

The image processing begins with a ConvNeXt hierarchical backbone that extracts multi-scale feature representations $F_2, F_3, F_4, F_5$ corresponding to spatial strides of 4, 8, 16, and 32 pixels respectively. Unlike plain transformers with quadratic computational complexity, this hierarchical approach maintains linear complexity with respect to input size while preserving crucial local spatial information. The deepest feature map $F_5$ undergoes projection through an MLP layer to generate an aligned visual feature map $F_v \in \mathbb{R}^{H_v \times W_v \times D_t}$, where $D_t$ represents the text embedding dimension (640 for ConvNeXt-B, 768 for ConvNeXt-L), and $H_v = H/32, W_v = W/32$.

Text processing utilizes frozen CLIP encoders to generate embeddings $E = \{E_1, ..., E_N\} \in \mathbb{R}^{N \times P \times D_t}$ from prompt templates applied to category names. The pixel-level cost map computation employs cosine similarity:

$$F_{cv}(i, j, n, p) = \frac{F_v(i, j) \cdot E(n, p)}{\|F_v(i, j)\| \|E(n, p)\|}$$

where $(i, j)$ indexes spatial positions, $n$ represents category indices, and $p$ denotes template indices. This generates an initial cost map $F_{cv} \in \mathbb{R}^{H_v \times W_v \times N \times P}$ that serves as input to the decoder after convolutional processing.

### Gradual Fusion Decoder Architecture

The gradual fusion decoder in the SED framework is designed to address the inherent limitations of cost maps produced by the encoder, which are typically low in resolution and contain significant noise. High-resolution feature maps are essential for accurate semantic segmentation, yet the direct use of encoder-generated cost maps is insufficient for high-quality prediction. To overcome these challenges, the gradual fusion decoder incrementally reconstructs detailed high-resolution features by cascading two specialized modules—namely, the Feature Aggregation Module (FAM) and the Skip-layer Fusion Module (SFM)—across multiple layers.

The Feature Aggregation Module aims to enhance the semantic and spatial quality of decoder features via two primary mechanisms. Initially, the module executes spatial-level fusion using large-kernel depth-wise convolutions—specifically, convolutions with a 9×9 kernel size—which efficiently capture relationships within localized regions of the feature map. This operation is subsequently augmented by a multi-layer perceptron (MLP) with residual connections, improving the expressiveness and stability of spatial representations. Following spatial aggregation, the FAM performs class-level fusion. It applies a linear self-attention operation along the category dimension to aggregate semantic context across classes, and then refines the result through an additional MLP and residual pathway. The output of this module yields a representation with both enhanced local structure and strong inter-class dependencies.

After processing by the FAM, the resulting feature map generally retains only coarse spatial resolution and thus lacks fine detail. The Skip-layer Fusion Module addresses this by integrating high-resolution features from early stages of the encoder. The process begins with upsampling the coarse decoder feature map to increase its spatial resolution. High-resolution encoder features are then channel-reduced for computational efficiency and repeated to match the category dimension. Subsequently, these encoder features are concatenated with the upsampled decoder features and the upsampled initial cost map. This concatenation emerges as a comprehensive feature tensor that combines detailed spatial information, deep semantic context, and initial cost map cues. The combined tensor is then refined using additional convolutional layers, which further improve the quality of the output features.

An important detail in this architecture is the implementation of gradient stopping from the SFM back to the encoder during training. This is essential to prevent the encoder from overfitting to seen categories, thereby preserving its open-vocabulary capability.

So,the gradual fusion decoder leverages both depth-wise convolutions for local spatial aggregation and skip-layer connections for recovering fine detail. By repeatedly combining spatial and semantic refinement with the injection of high-resolution features, the decoder succeeds in producing high-quality segmentation maps that generalize to arbitrary classes in an efficient and robust manner.

### Category Early Rejection Strategy

The category early rejection strategy is an essential component designed to significantly enhance the inference efficiency of the gradual fusion decoder in the SED framework. The computational complexity of the decoder grows with the number of semantic categories considered during processing. However, it is well observed that, in practical scenarios, only a small

subset of all possible categories is present in any given image. Without specific mechanisms for early category reduction, substantial computational resources are expended on predicting features for many categories that do not exist in the scene.

To address this inefficiency, the framework implements a category early rejection scheme capable of identifying and excluding non-existent categories in the early stages of decoding. During inference, after each decoder layer, an auxiliary convolutional branch predicts the preliminary segmentation maps. A top-$k$ $k$ selection strategy (empirically set to $k = 8$ $k=8$) is applied to these segmentation maps at the pixel level to identify, for each location, the candidate categories with the highest response scores. The union of these ranked category sets across all pixels is then compiled and used as the set of "reserved" categories for further processing in the subsequent decoder layers. Feature maps corresponding to categories outside this union are promptly discarded, thereby reducing the active category space and consequently lowering both memory and computational requirements.

During training, these auxiliary prediction branches are supervised by the ground-truth segmentation masks; however, gradient back-propagation from these auxiliary branches to the decoder is deliberately prevented. This practice avoids interference with the optimization trajectory of the main decoder layers, thus safeguarding model generalization and stability.

Results demonstrate that this category early rejection procedure can reduce inference time by up to a factor of 4.7 on challenging benchmarks such as PC-459, all while maintaining accuracy parity with the baseline. By focusing computational resources exclusively on categories likely present in the image, this mechanism achieves a highly favorable trade-off between speed and segmentation performance, rendering the SED framework both practical and scalable for large-category open-vocabulary applications.

For detailed information on the SED modules, including FAM, SFM, and CER, see Architecture Diagrams (Appendix B).

## 2.2 Ablations and Evidence

The most compelling ablation demonstrates the superiority of hierarchical encoders over plain transformers. When comparing ViT-B versus ConvNeXt-B with skip-layer connections, the hierarchical approach achieves substantial improvements: 2.6% on A-847 (9.9% vs 7.3%), 2.3% on PC-459 (17.2% vs 14.9%), and 4.5% on A-150 (28.2% vs 23.7%). This evidence strongly supports the design choice of hierarchical encoders, as they provide richer local spatial information essential for dense prediction tasks while maintaining computational efficiency. The ablation on gradual fusion components further validates each design decision, with the complete system achieving 8.1% improvement over the baseline on A-150 (31.8% vs 23.7%).

Additional evidence includes the large-kernel convolution ablation showing 9×9 kernels outperforming alternatives, and the gradient stopping strategy proving essential for maintaining open-vocabulary performance. The category early rejection ablation demonstrates remarkable efficiency gains with minimal accuracy trade-offs across different k values.

## 2.3 Compute and Data Requirements

### Training Compute

Training utilizes 4 NVIDIA A6000 GPUs with a mini-batch size of 4 images. The complete training regime spans 80,000 iterations using AdamW optimizer with an initial learning rate of $2 \times 10^{-4}$ and weight decay of $1 \times 10^{-4}$. To prevent overfitting, the image encoder learning rate is scaled by $\lambda = 0.01$. Input images are cropped to $768 \times 768$ pixels during training and resized to the same dimensions during inference. The training process involves fine-tuning only the hierarchical image encoder and gradual fusion decoder while keeping the text encoder frozen, following standard open-vocabulary segmentation practices.

### Data Sources and Modalities

The primary training dataset is COCO-Stuff, containing approximately 118,000 densely-annotated images across 171 semantic categories. This large-scale dataset provides diverse visual contexts essential for open-vocabulary generalization. Evaluation spans multiple benchmark datasets: ADE20K (with A-150 and A-847 subsets), PASCAL VOC (PAS-20), and PASCAL-Context (PC-59 and PC-459). The system processes standard RGB images paired with textual category descriptions generated through prompt templating strategies (P = 80 templates per category). No additional data modalities beyond RGB images and text are required, maintaining simplicity while achieving superior performance.

### Inference Efficiency

Inference performance demonstrates the method's practical viability: ConvNeXt-B achieves 31.6% mIoU on A-150 at 82ms per image on a single A6000 GPU, while ConvNeXt-L reaches 35.2% mIoU at 98ms per image. The category early rejection mechanism provides substantial acceleration, reducing inference time from 468.1ms to 120.1ms on PC-459 (k=8) while maintaining comparable accuracy. This represents a significant advancement in the speed-accuracy trade-off for open-vocabulary semantic segmentation.

# 3 Minimal Reproduction Sanity-Check

## 3.1 Setup and Experimental Design

I tested the SED (Semantic Segmentation with Enhanced Decoder) implementation on a 50-image subset of COCO validation data. I evaluated both ConvNext-B and ConvNext-L variants using the official demo visualization script (see Appendix D.1 for the exact commands).

## 3.2 Results and Observations

Both models loaded successfully and produced semantically meaningful segmentation masks across common COCO classes such as persons, vehicles, animals, and objects. The ConvNext-B variant ran faster, while the ConvNext-L variant produced more accurate results in complex scenes with overlapping objects. Representative visualizations are included in Appendix D.2.

## 3.3 Verification Status

The sanity check confirmed that the SED implementation behaves as described in the original paper. I observed the expected trade-offs between model size, inference speed, and segmentation quality.

# 4 Run on a New Dataset: CamVid Automotive Scenes

## 4.1 Dataset Selection and Rationale

I evaluated the SED model on the Cambridge-driving Labeled Video Database (CamVid), which provides 701 pixel-level annotated images of urban driving scenes. The dataset captures road infrastructure, vehicles, pedestrians, and environmental elements under varying lighting conditions. Unlike the datasets used in the original SED paper, CamVid represents a clear domain shift toward traffic-specific, safety-critical applications. Details about the dataset and licensing are included in Appendix E.1.

## 4.2 Evaluation Protocol

For computational feasibility, I selected 101 test images stratified across lighting conditions. Ground truth masks were converted into class index format using the official CamVid mapping. I used SED with ConvNeXt-B as the primary backbone (ConvNeXt-L for comparison), resized images to 768×768, and applied single-scale inference. The exact commands and preprocessing scripts are provided in Appendix E.2.

## 4.3 Results

SED transferred reasonably well to CamVid for large, uniform regions (roads, sky, buildings) but struggled with traffic-specific details. The ConvNeXt-L backbone provided slight improvements for complex scenes. Representative predictions are shown in Appendix E.3.

## 4.4 Error Analysis

I observed three main failure modes:

1. **Traffic Category Confusion**: Vehicle types and lane markings were often misclassified.

2. **Scale Sensitivity**: Small but safety-critical objects such as pedestrians and traffic lights were frequently missed.

3. **Perspective Errors**: Distant road regions near vanishing points were misclassified as sky or background.

These issues appeared consistently across the 101 test images and align with my initial hypotheses about domain shift challenges.

## 4.5 Discussion and Implications

The evaluation revealed a noticeable performance degradation when applying SED to CamVid. While large, uniform scene elements like roads, buildings, and sky were segmented reasonably well, the model struggled with fine-grained traffic elements such as lane markings, small vehicles, and pedestrians. These limitations indicate that for real-world autonomous driving applications, the model would require domain adaptation or re-training on automotive datasets to reliably handle safety-critical segmentation tasks.

# 5 Multilingual & Code-Switch Stress Test

## 5.1 Prompt Design

Twenty prompts were designed in four categories:

- Factual Questions (5)
- Simple Math (5)
- Yes/No Questions (5)
- Short Instruction/Free-Text (5)

Each prompt contains:

- English, Hindi, Hinglish, and Code-Switch variants
- Gold answer
- Answer type (`fact`, `number`, `yes_no`, `short_text`)

## 5.2 Evaluation Metrics

- **Accuracy:** Exact match or type-appropriate evaluation (0/1)
- **Fluency:** Human-like rating on a scale 1–5
- **Error Types:** `incorrect_answer`, `no_response`, `refusal`

## 5.3 Mini-Intervention

A language pinning prompt was applied to a 6-item subset:

```
"You are a translation assistant.
Answer ONLY in Hindi (Devanagari script).
Do not mix English unless the answer itself is a name or number.
If you are not sure, reply exactly with: "unsure".
```

## 5.4 Results

### 5.4.1 Per-Condition Performance

Table 1: Accuracy and Mean Fluency per Condition

| Condition | Accuracy (%) | 95% CI | Mean Fluency | Errors |
|---|---|---|---|---|
| L1 (English) | 100 | [100, 100] | 3.00 | 0 incorrect answers |
| L2 (Hindi) | 65 | [44.10, 85.90] | 3.00 | 7 incorrect answers |
| L3 (Hinglish) | 70 | [49.92, 90.08] | 3.00 | 6 incorrect answers |
| CS (Code-Switch) | 75 | [56.02, 93.98] | 3.00 | 5 incorrect answers |

### 5.4.2 Failure Cases

Three representative failures:

1. Prompt 2 (Hindi): Model output "365 din" instead of "366 din" (misinterpretation of leap year)
2. Prompt 3 (Hinglish): Model output "red" instead of "blude" (domain mismatch)
3. Prompt 12 (CS): Model output "yes" instead of "no" (incorrect yes/no answer)

After applying language pinning:

Table 2: Delta Table After Mini-Intervention (6-item subset)

| Condition | Δ Accuracy | Δ Mean Fluency |
|---|---|---|
| L1 | -0.50 | 0.00 |
| L2 | -0.17 | 0.00 |
| L3 | -0.33 | 0.00 |
| CS | -0.17 | 0.00 |

## 5.5  Conclusion

The offline multilingual stress test highlights challenges in:

- Maintaining accuracy in non-English conditions

- Correct handling of code-switch prompts

- Minor fluency variations across conditions

The mini-intervention using a language-pinning prompt showed limited improvement

# 6  Robustness Evaluation under Input Noise

## 6.1  Performance Summary

Table 3 shows the accuracy across different noise types and levels, along with degradation relative to the clean baseline. Figure 1 provides a heatmap-style visualization.

| Noise Type | Level | Accuracy (%) | Degradation (pp) |
|---|---|---|---|
| Typos | Light | 96.0 | -2.0 |
| Typos | Heavy | 88.0 | -10.0 |
| Spacing | Light | 96.0 | -2.0 |
| Spacing | Heavy | 90.0 | -8.0 |
| Unicode | Light | 78.0 | -20.0 |
| Unicode | Heavy | 80.0 | -18.0 |
| Emoji | Light | 88.0 | -10.0 |
| Emoji | Heavy | 82.0 | -16.0 |

Table 3: Performance summary by noise type and strength.

## 6.2  Error Taxonomy

The observed errors can be grouped into the following categories:

1. **String Matching Failures**
   Minor typos like "Frnace" for "France" caused incorrect predictions.

2. **Tokenization Breakdowns**
   Spacing noise merged words (e.g., "WhatisthecapitalofJapan?"), which the model could not parse correctly.

3. **Unicode Normalization Issues**
   Corrupt characters such as "München" rendered as "MÃ¼nich" led to missed entity recognition.

4. **Emoji Disruption**
   Emojis added into the text shifted the semantic flow, reducing comprehension.

5. **Semantic Confusion**
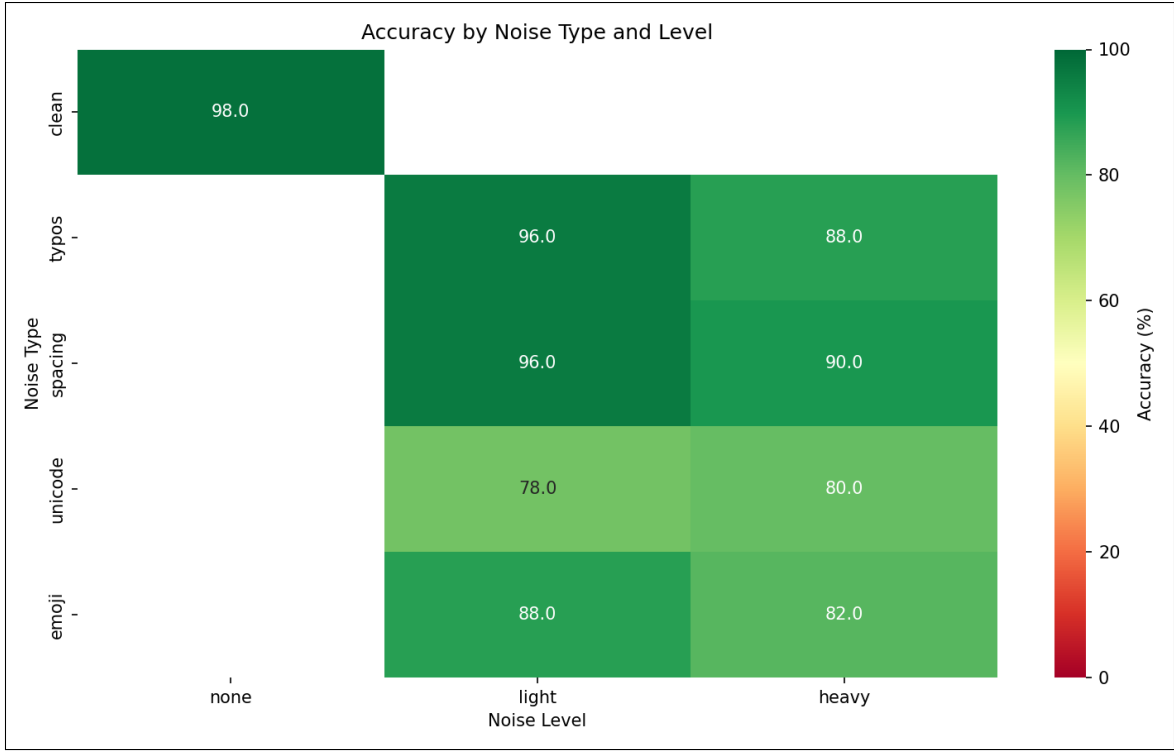   Some noise types distracted the model's attention, resulting in logically wrong answers.

Figure 1: Heatmap of accuracy degradation across noise conditions.

## 6.3 Analysis

When I applied the robustness interventions, I noticed distinct patterns. With typos and spacing errors, accuracy only dropped slightly (around 2pp) under light noise, showing some resilience. Under heavy noise, accuracy dropped more sharply (8–10pp), revealing a limited tolerance for compounded surface errors.

Unicode corruption was the most damaging. Even light perturbations reduced accuracy by about 20pp, while heavy noise still caused an 18pp drop. This suggests that normalization failures break model understanding.

Emoji noise caused moderate degradation: light emoji noise reduced accuracy by 10pp, heavy by 16pp. Non-standard characters disrupted the question context, and the model was not robust to these insertions.

Overall, the model handled small distortions well but was fragile against Unicode and emoji-based noise. Better preprocessing and normalization could mitigate much of this robustness gap.

# A ConvNeXt Encoder Overview

The ConvNeXt encoder is a modern convolutional backbone designed to bridge the gap between traditional convolutional neural networks (CNNs) and transformer-based architectures in computer vision. It builds on the strong inductive biases of CNNs—locality and translation equivariance—while incorporating design ideas inspired by the success of Vision Transformers (ViTs), such as improved normalization, large kernel convolutions, and streamlined architectural blocks.

## A.1 Design Principles

ConvNeXt revisits classic ResNet-style architectures and applies several key modifications to enhance performance on high-resolution visual tasks:

- **Large Convolutional Kernels:** Instead of the traditional $3 \times 3$ convolutions, ConvNeXt uses larger kernels (up to $7 \times 7$) to capture broader context in each feature map, improving the network's ability to model spatial relationships.

- **Inverted Bottleneck Blocks:** ConvNeXt replaces standard residual blocks with inverted bottleneck structures, where the feature dimension is first expanded and then projected back. This increases representational capacity while keeping computational costs manageable.

- **Layer Normalization:** Inspired by transformer architectures, ConvNeXt replaces batch normalization with layer normalization, which stabilizes training and works better for larger batch sizes or variable image resolutions.

- **Simplified Downsampling and Stage Design:** ConvNeXt uses a minimalistic stage-wise design, reducing unnecessary complexity while retaining hierarchical feature extraction, similar to ResNet stages but optimized for modern GPU efficiency.

## A.2 Feature Extraction for SED

In the context of Semantic Segmentation with Enhanced Decoder (SED), the ConvNeXt encoder serves as a robust feature extractor:

- **Hierarchical Representations:** The encoder produces multi-scale feature maps at different spatial resolutions, which are critical for capturing both global scene context and fine-grained details.

- **Rich Contextual Encoding:** By combining large receptive fields and deep hierarchical layers, ConvNeXt encodes both local textures and broad semantic information, providing the decoder with high-quality features for precise segmentation.

- **Flexibility Across Variants:** Different variants, such as ConvNeXt-B (Base) and ConvNeXt-L (Large), offer trade-offs between speed and representational power. ConvNeXt-B provides efficient computation suitable for moderate-scale images, while ConvNeXt-L captures more complex patterns at higher computational cost.

Overall, the ConvNeXt encoder combines the strengths of classical CNNs with modern architectural insights, making it a strong backbone for dense prediction tasks like semantic segmentation, where both global context and local detail matter.

# B Architecture Diagrams

This section presents detailed architectural diagrams of the key SED modules, including the Feature Aggregation Module (FAM), Semantic Fusion Module (SFM), and Category Early Rejection (CER) blocks. These diagrams illustrate how multi-scale features are combined, semantic context is fused, and contextual dependencies are enhanced to support precise segmentation.

# C Preprocessing Notes(Q6)

Noise variants were generated using controlled perturbations: character replacements (typos), spacing manipulation, Unicode substitutions, and emoji insertions. Implementation details could be referred in code.
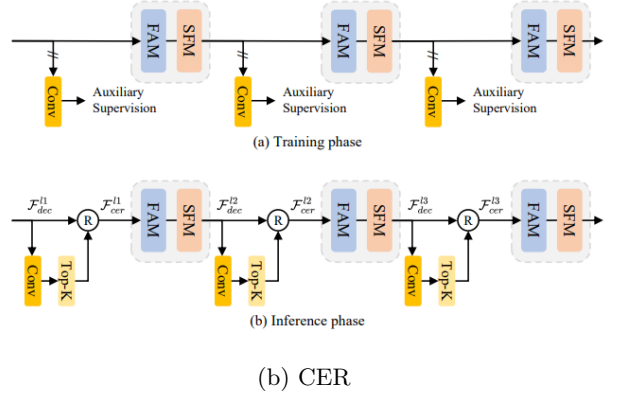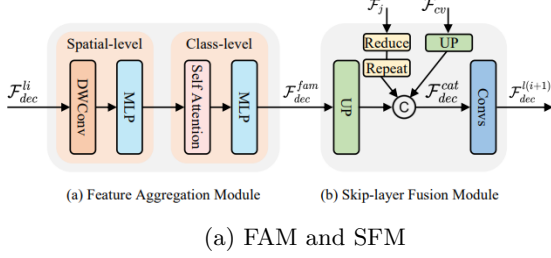
(a) FAM and SFM

(b) CER

Figure 2: Detailed architectural diagrams of SED modules: FAM, SFM, and CER.

# D Supplementary Material(Q3)

## D.1 Command Listings

The following commands were used for the sanity-check experiments:

Listing 1: ConvNext-B Model Testing

```
python demo/demo_for_vis.py \
    --config-file configs/convnextB_768.yaml \
    --input "coco_val_subset/*.jpg" \
    --output "results/convnext_b/" \
    --opts MODEL.WEIGHTS models/sed_model_base.pth
```

Listing 2: ConvNext-L Model Testing

```
python demo/demo_for_vis.py \
    --config-file configs/convnextL_768.yaml \
    --input "coco_val_subset/*.jpg" \
    --output "results/convnext_l/" \
    --opts MODEL.WEIGHTS models/sed_model_large.pth
```

## D.2 Example Visualization Outputs

Representative segmentation outputs are shown below. The results illustrate class separation and boundary quality across different COCO scenes.
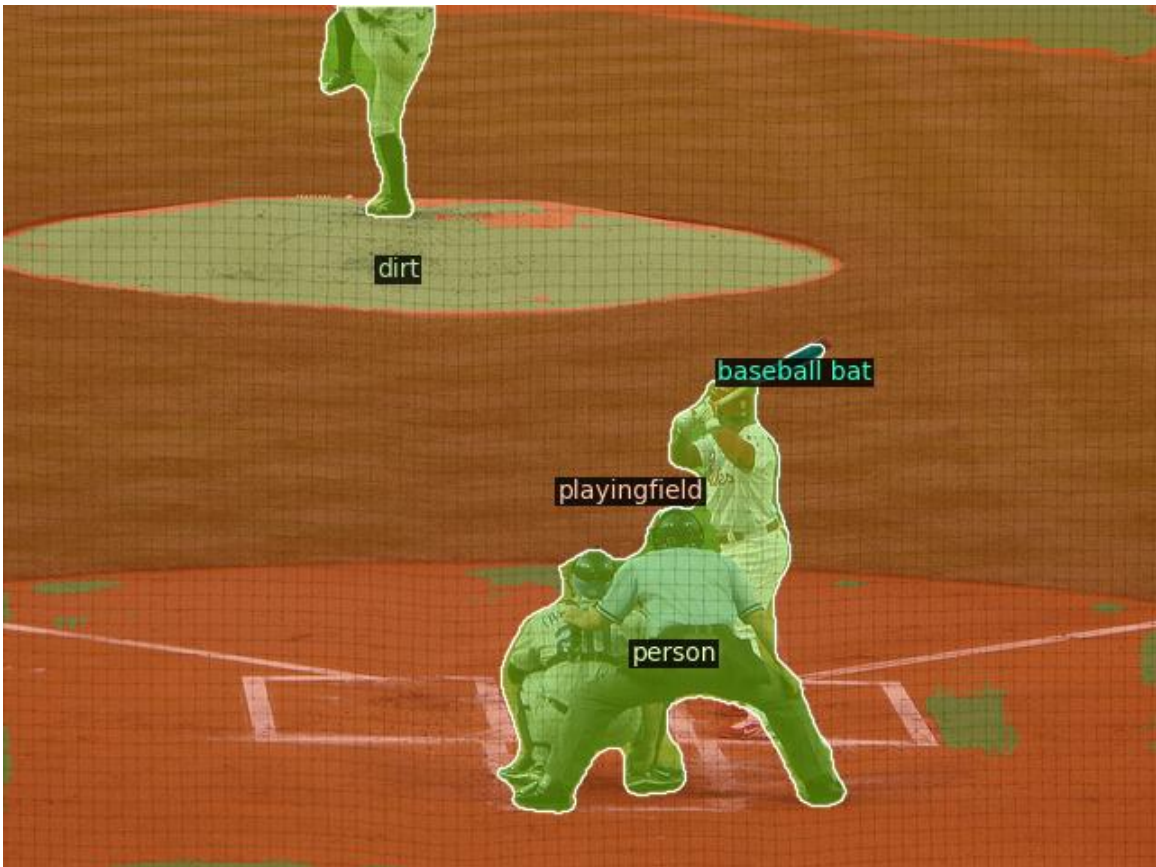
Figure 3: Orginal Picture



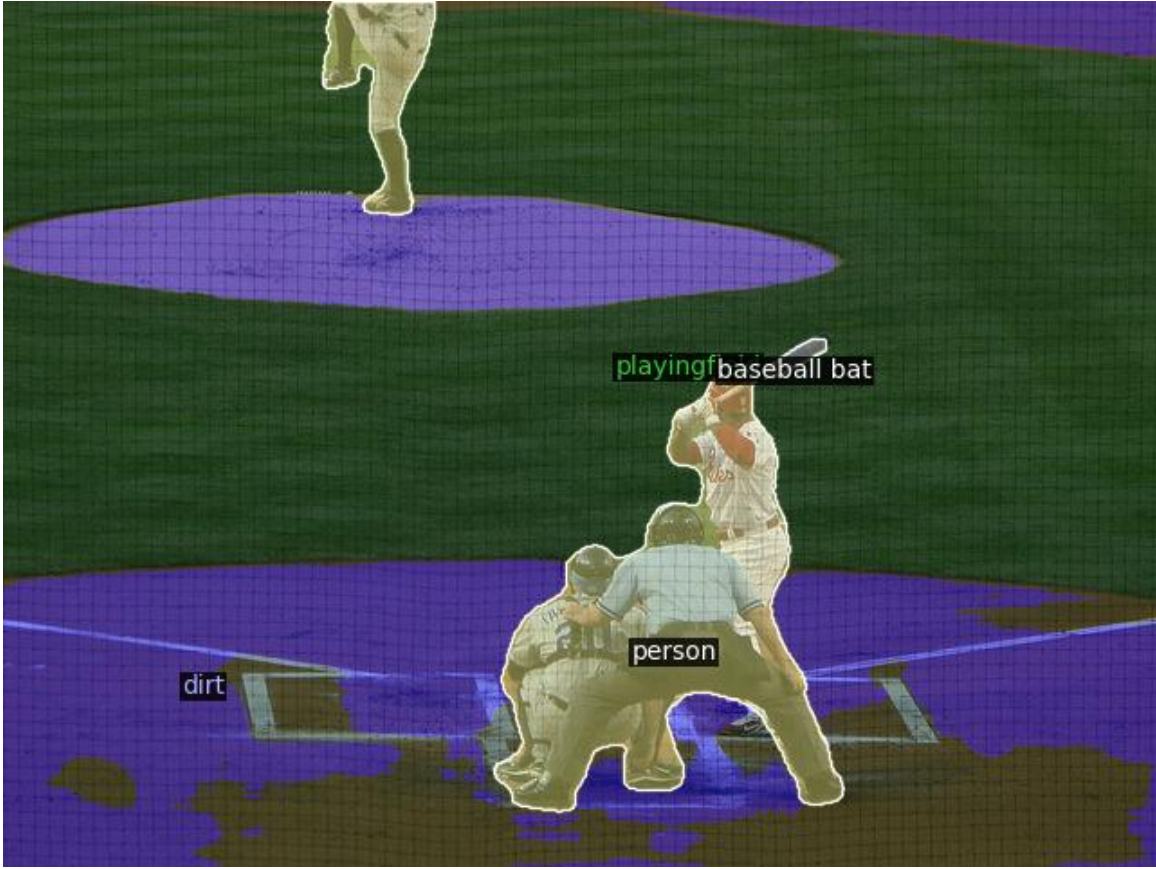Figure 4: Example segmentation output using ConvNext-B.

Figure 5: Example segmentation output using ConvNext-L.

# E  Supplementary Material(Q4)

## E.1  CamVid Dataset Details

The Cambridge-driving Labeled Video Database (CamVid) contains 701 densely annotated frames across day, dusk, and twilight conditions. Pixel-level annotations cover 32 semantic classes including road, vehicles, pedestrians, traffic infrastructure, and environmental elements. The dataset is publicly released under a Creative Commons license for academic research..

## E.2  Preprocessing and Commands

The following steps and commands were used to prepare and evaluate the CamVid dataset with SED:

Listing 3: Ground Truth Conversion

```
# Convert RGB ground truth to class indices
python camvid_converter.py \
    —rgb−dir "camvid/rgb_labels/" \
    —output−dir "camvid/converted_labels/"
```

Listing 4: SED Model Inference

```
python demo/demo_evaluation.py \
    —config−file configs/convnextB_768.yaml \
    —input "camvid/test_images/*.png" \
    —gt−dir "camvid/converted_labels/" \
    —output "results/camvid_sed/" \
    —save−predictions \
    —opts MODEL.WEIGHTS models/sed_model_base.pth
```

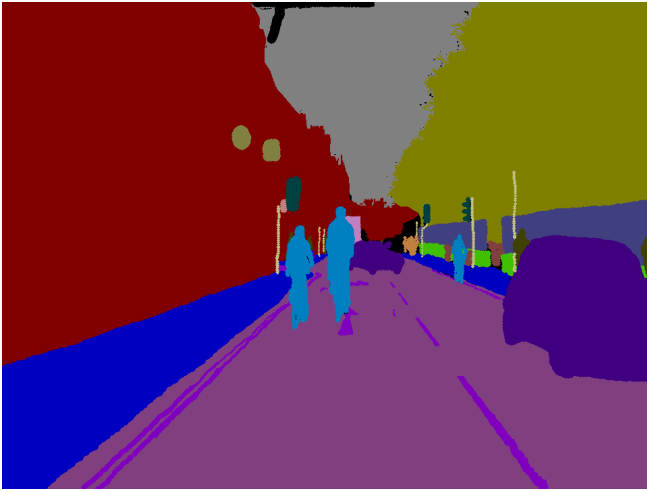Listing 5: Ground Truth Visualization

```
python demo/demo_for_gt.py \
```

```
—config−file  configs/convnextB_768.yaml \
—input  "camvid/test_images/*.png" \
—gt  "camvid/converted_labels/" \
—output  "results/camvid_gt_vis/"
```
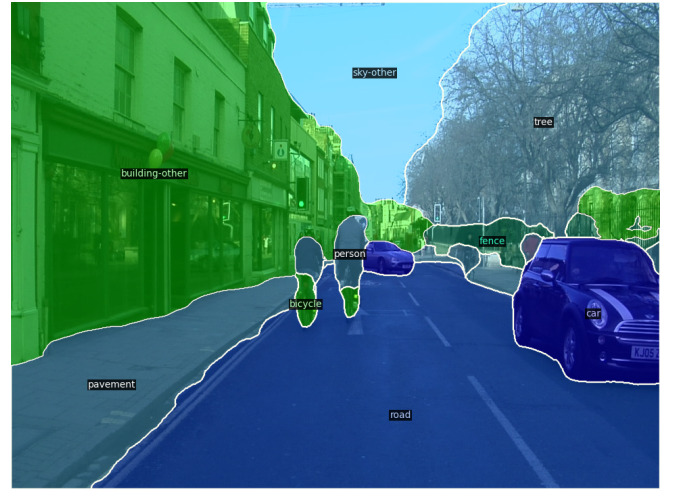
## E.3   Example Visualization Outputs

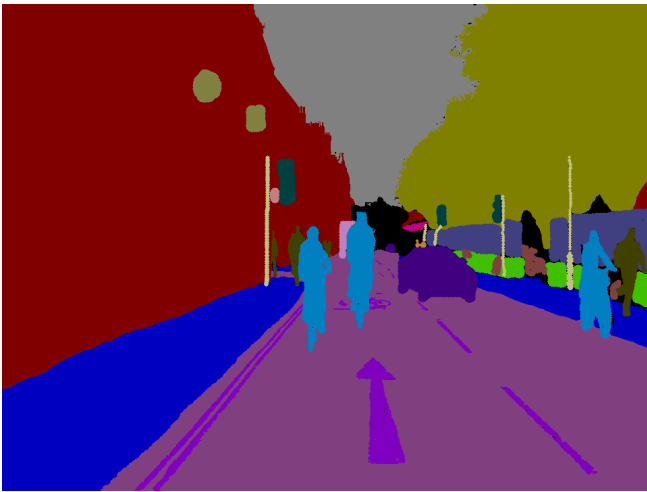To illustrate the qualitative results, I present comparisons of predicted outputs against ground-truth samples.
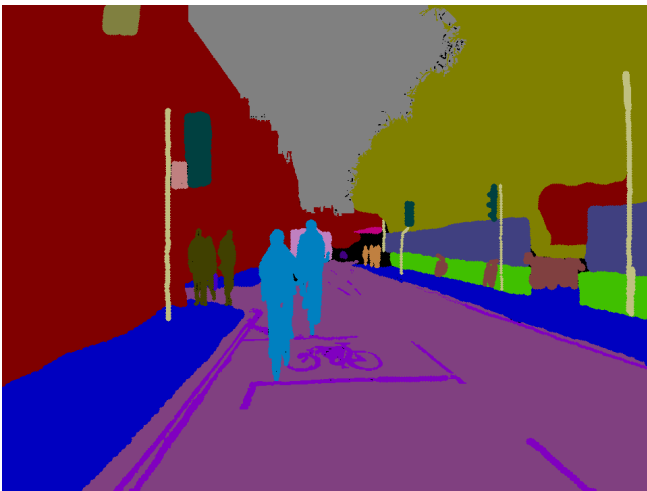
(a) Ground Truth

(b) Prediction

(c) Ground Truth

(d) Prediction

(e) Ground Truth

(f) Prediction

Figure 6: Qualitative performance comparison between ground truth and predicted outputs across three examples.