The rising STAR of Texas

UNIVERSITY

TEXAS

Cascaded Workflow for Automated Pavement Crack Detection and Measurement Using Generative **AI Models and Image Processing Techniques**

Abstract

In this research, Generative AI (Artificial Intelligence) models such as Ground DINO and Florence 2 are developed in the mission of automated pavement distress detection with 2D/3D image data, and computer vision techniques are applied for the quantification of pavement cracks. In this cascaded workflow, the Generative AI models are compared with different pavement distress detection models and developed with the dataset collected and annotated by the research team. After the surface defects are detected, a technique using the Structured Forest edge extraction method is proposed to delineate the shape of the cracks and then measure the crack geo-information (widths, lengths, and area) at the distress locations. Finally, weighted cracking widths, lengths, and areas are applied to calculate the distress severity levels or distress scores to follow the protocols of distress identification guidelines and pavement management standards specified by the U.S. and state DOTs. This two-step workflow with the Generative AI and image processing techniques will provide an automated pipeline to process the field collected pavement image data and assess the cracking condition of the pavement surface with costeffectiveness. Introduction **1. Traditional and open-vocabulary object detection** Text ≻ Vocabulary Online Vocabulary Encoder **▲ ▲ Object Detector** Large Detector Figure 1 Fixed-vocabulary (a) and open-vocabulary (b) detection 2. Models used for pavement distress detection 1) Real-Time Detection Transformer (RT-DETR) Images are the input. Intrascale feature interaction (AIFI) and crossscale feature-fusion module (CCFM) are the encoder. IoU-aware query is used for image feature selection. The decoder generating Fransverse boxes and confidence scores. 2) YOLOv12 YOLOv12 is an attention-centric variant of the YOLO (You Only Look Once) family. 3) YOLO World Text and the input image are encoded with YOLO backbone. 4) Grounding-DINO With ground pre-training, the Grounding DINO (Distillation with No Labels) can utilize a feature enhancer, a language-guided query selection and a modality fusion with a cross-modality decoder. 5) Florence 2: The semantic granularity at three stages forms a comprehensive annotation module named FLD-5B. Florence-2 is a unified architecture based on FLD-5B. Semantic Granularity tycle on a road with a red car in the kground. The road is lined with ees on both sides and there is another erson riding another bicycle in front of ter. The date "9/22/2023" is visible in visual grounding ansverse crack. Longitudinal crack. detailed caption & object detection FLD-5B Comprehensive Annotation nt, sealed transverse. Sea riding a bil down a street next ngitudinal, Lane Longitudinal to a red ca person Florence-2 car & object detection (Unified Architecture)

Figure 2. An illustration of Florence 2 enabling semantic granularity and spatial hierarchy as a vision foundation model.

Region-level

mage-level

Pixel-level Hierarchy

Figure 7 Generative AI models (vision-language models) for pavement distress detection.

Kev and values

Yongsheng Bail and Feng Wang Ingram School of Engineering, Texas State University



	Without tiling			With tiling		
Class name	Precision	Recall	mAP50	Precision	Recall	mAP50
Transverse crack	0.344	0.482	0.354	0.663	0.602	0.632
Sealed transverse crack	0.517	0.254	0.296	0.567	0.516	0.53
Joint	0.204	0.669	0.421	0.82	0.815	0.847
Longitudinal crack	0.325	0.253	0.197	0.539	0.481	0.459
Sealed longitudinal crack	0.257	0.349	0.245	0.657	0.573	0.589
Lane longitudinal crack	0.258	0.289	0.145	0.588	0.551	0.518
Block crack	0.0486	0.0339	0.0184	0.371	0.285	0.249
Alligator crack	0.0991	0.0118	0.0603	0.614	0.509	0.538
Failures	0.376	0.268	0.289	0.621	0.651	0.663
Overall performance	0.27	0.29	0.225	0.604	0.553	0.558

	3D only			Fused 2D/3D		
Class name	Precision	Recall	mAP50	Precision	Recall	mAP50
Transverse crack	0.647	0.615	0.641	0.663	0.602	0.632
Sealed transverse crack	0.568	0.54	0.506	0.567	0.516	0.53
Joint	0.824	0.821	0.859	0.82	0.815	0.847
Longitudinal crack	0.544	0.488	0.485	0.539	0.481	0.459
Sealed longitudinal crack	0.622	0.541	0.582	0.657	0.573	0.589
Lane longitudinal crack	0.556	0.514	0.478	0.588	0.551	0.518
Block crack	0.355	0.259	0.225	0.371	0.285	0.249
Alligator crack	0.584	0.463	0.507	0.614	0.509	0.538
Failures	0.514	0.687	0.624	0.621	0.651	0.663
Overall performance	0.579	0.548	0.545	0.604	0.553	0.558



del name	3D or fused	Data augmentation	mAP50	Precision	Recall		
DETR	3D	yes	0.545	0.579	0.548		
_DETR	fused	yes	0.558	0.604	0.553		
LOv12	fused	yes	0.546	0.610	0.526		
LO World	fused	yes	0.511	0.570	0.516		
unding-DINO	fused	yes	0.515	0.522	0.504		
rence 2	3D	yes	0.301	0.400	0.440		
rence 2	fused	yes	0.263	0.450	0.360		

This research is part of the project "using artificial intelligence to improve the accuracy of automated pavement condition data collection", funded by the U.S. National Science Foundation (grant No. 2213694).

This research is also sponsored by Texas Department of Transportation (project No. 0-7150).