









Computer Vision for Digital Twins (CV4DT), where perception meets geometry

Dr Olaf Wysocki Assistant Research Professor

CV4DT
Department of Engineering
University of Cambridge









About me

- Assistant Research Professor (08/2025 now)
 - @ University of Cambridge, leading the CV4DT group
- PostDoc + PhD (09/2020 08/2025)
 - @ Technical University of Munich
 - @ Photogrammetry & Remote Sensing, Prof Uwe Stilla 3D Semantic Understanding / 3D Object Reconstruction
- Previous stints in industry (e.g., Audi)



















Dr Yixiong Jing Postdoctoral Research Associate



Dr Brian Sheil Laing O'Rourke Associate Professor in Construction Engineering



Postdoctoral Research Associate ≥ g ()

Researchers



Haibing Wu PhD Candidate **2** 2



Dr Olaf Wysocki Group Lead, Assistant Research Professor







2 2



Wanru Yang PhD Candidate









Daniel Lehmberg Intern





Ziyang Xu PhD Candidate



Qizhen Ying PhD Candidate





Qilin Zhang PhD Candidate









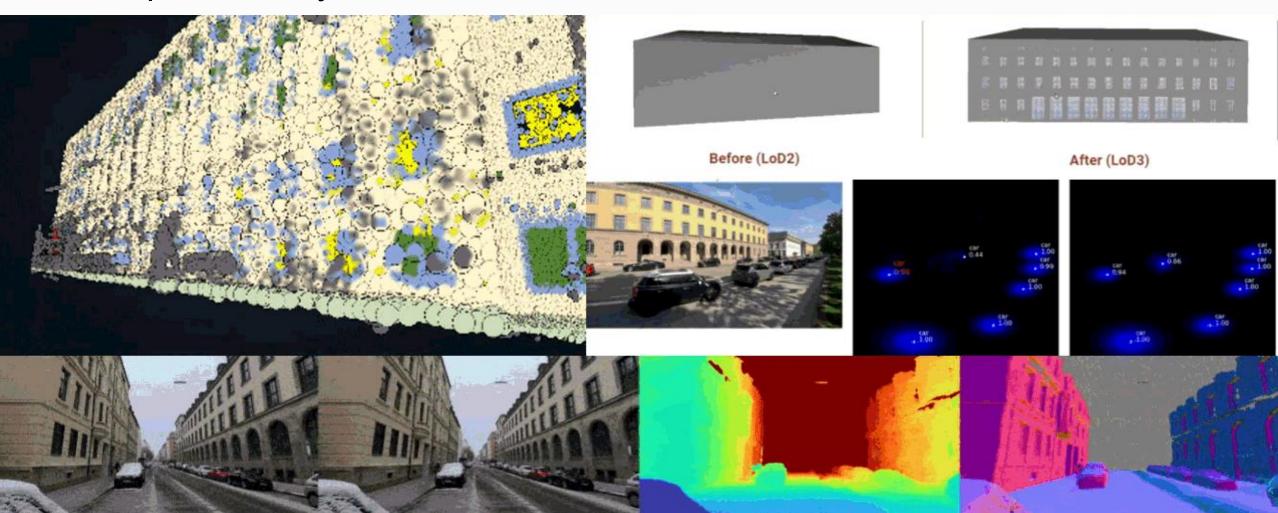








Perception is all you need





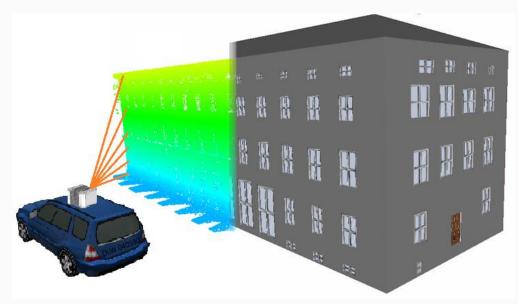


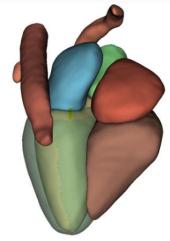




Motivation

- Digital Twins (DTs) necessary, e.g.,:
 - Gaming
 - Medical
 - Garment
 - Up to country-wide city models
- But DTs require:
 - 3D minimum-viable representation (think CAD)
 - Semantics
 - Accuracy & Completness
- 3D Computer Vision so-far:
 - Small-scale and...
 - Toy examples











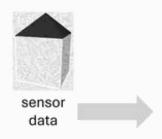


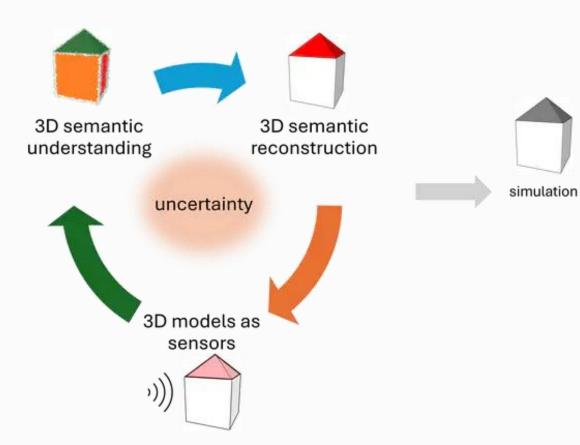


CV4DT Vision

Focal points:

- 3D semantic understanding,
- 3D semantic reconstruction,
- 3D models as sensors,
- Uncertainty quantification overarching all three aspects.













3D Computer Vision for Built Environment

Can we go beyond the toy examples?



a plush dragon toy



Construction site



3D model (mesh)



3D model (BIM)









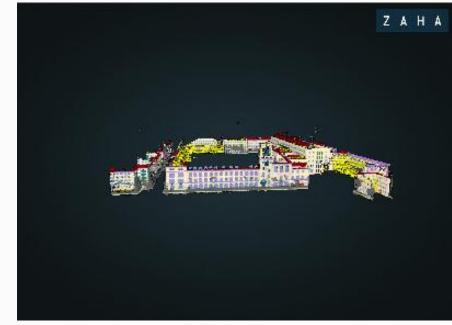
Benchmark Datasets: ZAHA

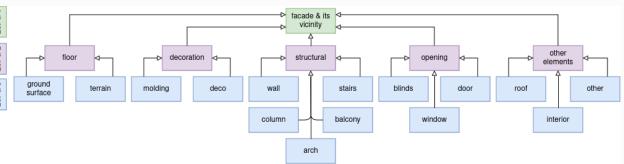
ZAHA: The largest facade segmentation dataset

Level of Facade Generalization (LoFG)

- 15 façade classes
- 3 levels of generalization
- Over 600 mln points









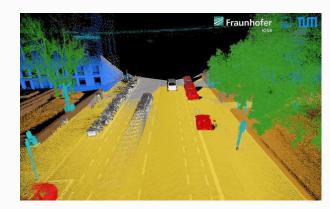






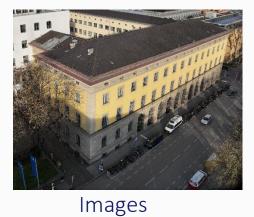
Benchmark Datasets: TUM2TWIN https://tum2t.win

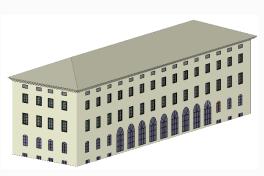


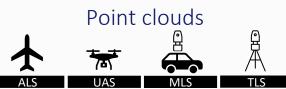


















3D models











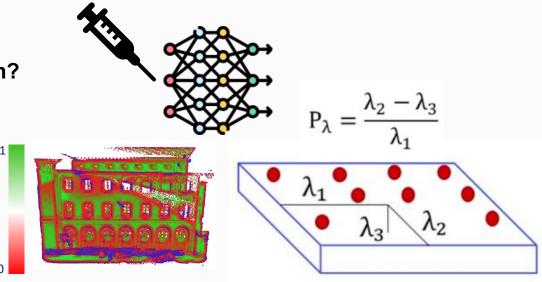




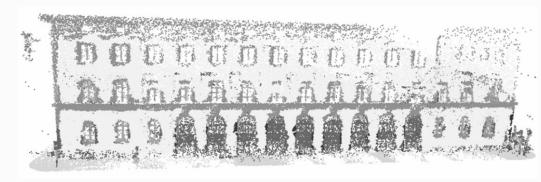
3D Semantic Segmentation

EARLy: How effectively can we use geometry for segmentation?

- Fusion of geometric features + learned features
 - EARLy + PointNet: ~81% (+9%)
 - EARLy + PointNet++: ~81% (+6%)
- EARLy + Point Transformer (PT): ~88% (+10%)
- Faster convergence (2x Faster)
- Less compute needed (2x Less Parameters)



Encoding planarity



Input - point cloud



Output - segmented point cloud





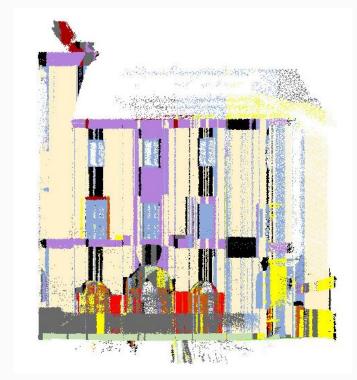




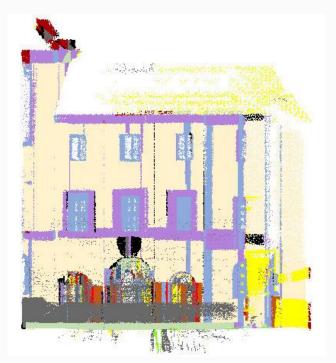
3D Semantic Segmentation

EARLy: How effectively can we use geometric prior?

Example of façade segmentation results



Before our EARLy



After our EARLy





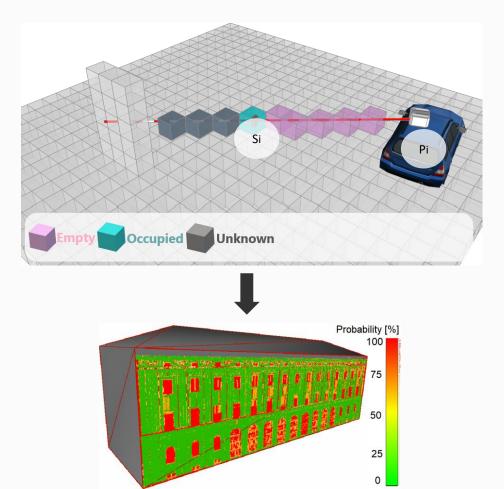




Laser scanner rays penetrate transparent objects -> no point cloud!

Can we use it to our advantage?

- Probabilistic ray tracing of points -> Empty, Occupied, Unknown voxels
- Comparison to 3D building model -> Confirmed, Conflicted voxels
- Conflicted voxels -> refinement required







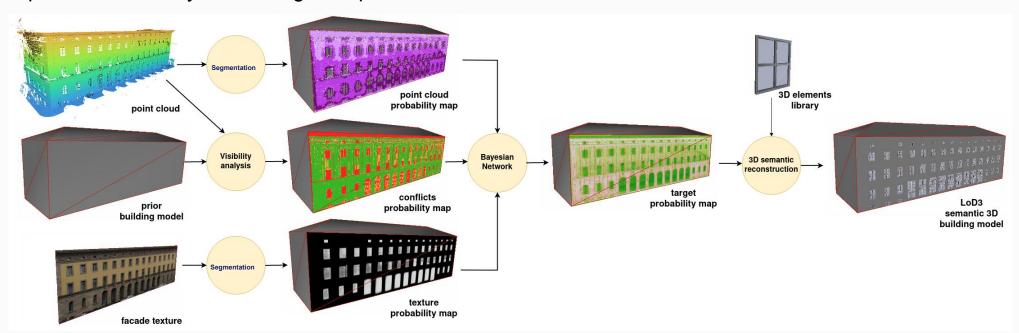




How effectively can we use our strong prior for segmentation and reconstruction?

Scan2LoD3: Uncertainty-aware late fusion of semantic maps for more accurate 3D reconstruction

- Bayesian Network calculates target based on uncertain input
- Output -> Uncertainty-aware target map and reconstruction instances











Scan2LoD3:

Conflict Map + Point Cloud Map:

Openings IoU: ~61% (+54% / +3%)

Conflict Map + Point Cloud Map + Texture Map:

ground-truth

Openings IoU: ~76% (+59% / +18%)

Point-cloud-only:

Openings IoU: 7%

Image-only:

Openings IoU: 58%



windows with blinds IoU = 73.8IoU = 51.1w/images

w/o images

IoU = 46

IoU = 53.6





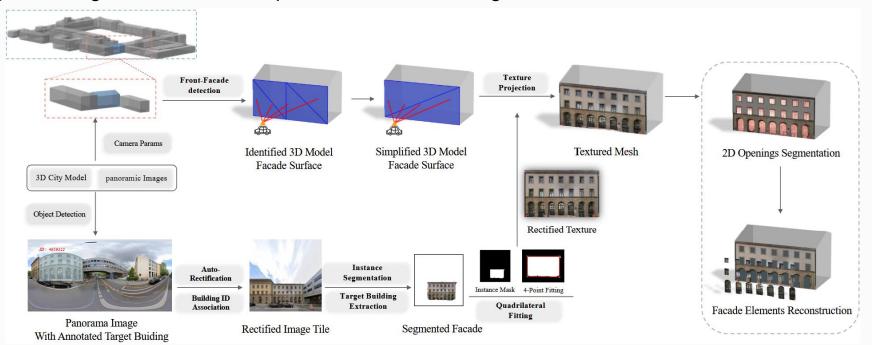




What happens if we do not have any scanners? Let use panoramic images (think Google Street View)!

Texture2LoD3: More accurate panoramic image projection -> more accurate segmentation & reconstruction

- Building models composed of planes use them as rectification targets for panoramic images
 - Image2building and transform to quasi-orthorectified image





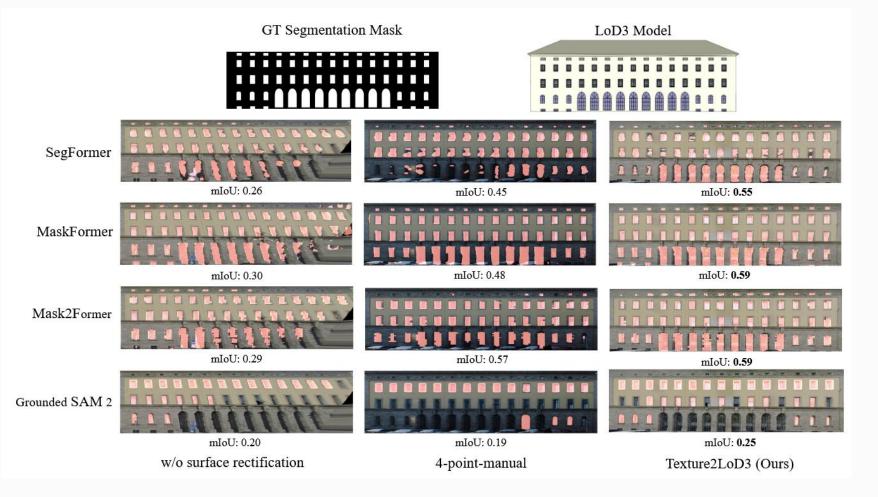






Texture2LoD3:

- +10% IoU vs baselines
- On par with manual













Gaussian Splatting reconstruction yielding appealing visualization



Gaussian Splatting reconstruction in the practical, built environment scenario









3D Models as Sensors & Reconstruction

Gaussian Splatting (GS) – State-of-the-art









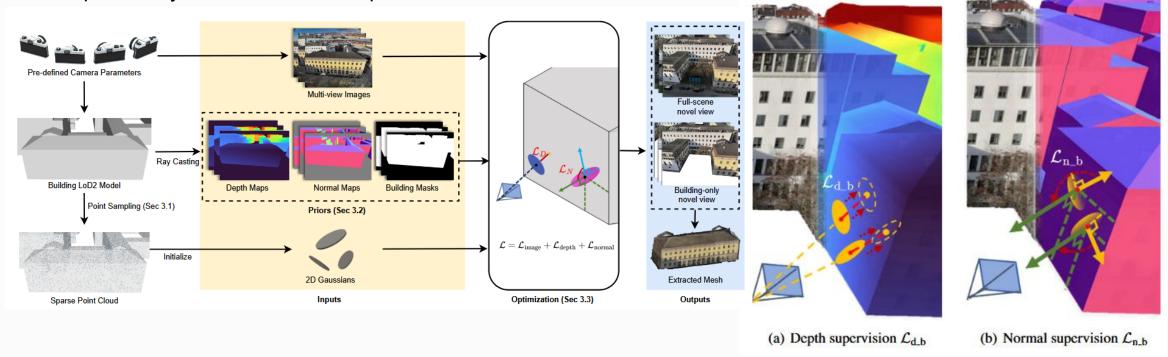


3D Semantic Reconstruction & 3D Models as Sensors

Prior-driven Gaussian Splatting

GS4B: Can we recycle worldwide-available 3D models to initialize + optimize GS reconstruction?

- Leverage prior model approximations to reduce noise and outliers
- Supervise by model-extracted depth and normals









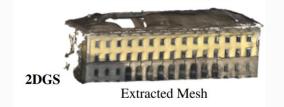


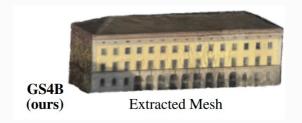
3D Models as Sensors & Reconstruction

GS4Buildings:

- +33% higher accuracy,
- +64% higher completness













Our GS4B (Building-only mode)







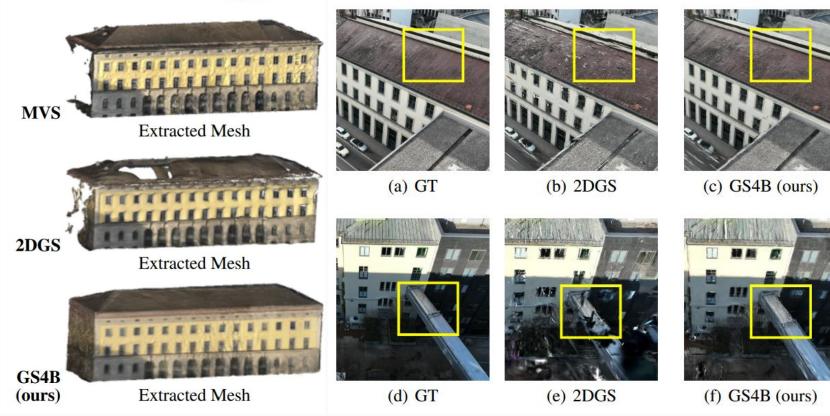


3D Semantic Reconstruction & 3D Models as Sensors

GS4B:

- +33% higher accuracy,
- +64% higher completness













3D Uncertainty-aware Reconstruction

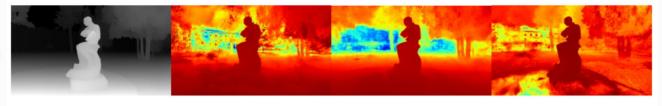
CD-GS: Uncertainty-aware Gaussian Splatting Including confidence for RGB edge + depth maps

Confidence-Aware Depth Regularization:

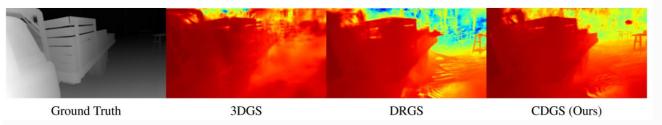
- Generates confidence maps for each depth map
- Adaptively adjusts depth loss during optimization
- Enables reconstructing:
 - Small objects
 - Corner cases (literally as well!)



















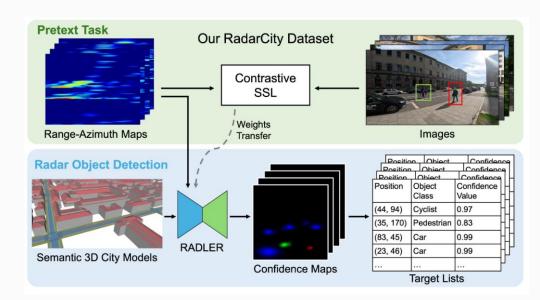
3D Models as Sensors

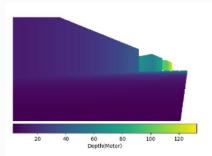
Can we use 3D semantic models for radar object detection?

Yes!

It turns out that:

- We get the depth estimation + semantics of the scene for free
- SSL used for getting the dynamic objects features
- We can improve detection results by ca. +7% vs SOTA





(a) Depth map includes the distances to objects in the scene.



(c) Camera image.



(b) Semantic map represents scene segmentation, where each color encodes a distinct object class.



(d) Semantic map overlaid on the image.









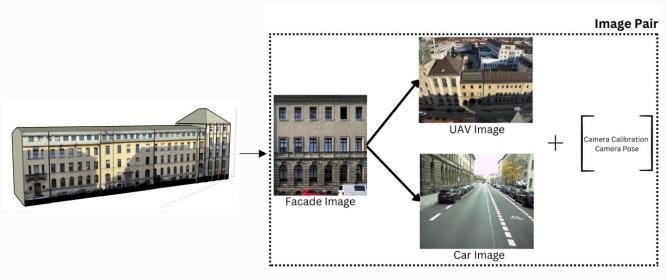
3D Models as Sensors

Can we use 3D semantic models for localization without GNSS? Yes!

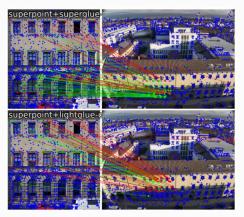


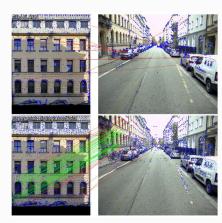
2D image to 3D geometry problem

- We get the 3D global (worldwide) coordinates for free
- Drastic changes in representation still challenging
- Still possible to retrieve inliers (+optional radiometry in limited cases)















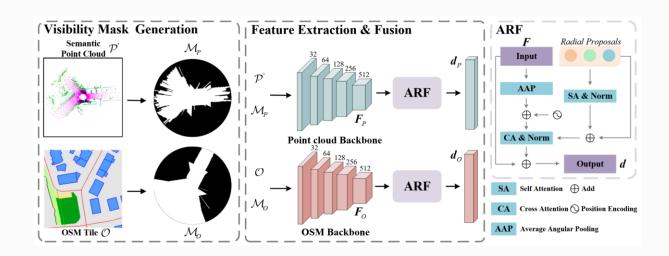


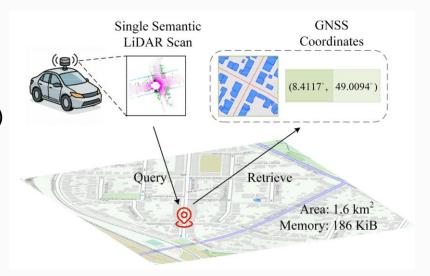
3D Models as Sensors

Can we use semantics of OpenStreetMap for LiDAR place recognition?

Turns out that we need only one scan to get meter-level accuracy

- Generate visibility mask to attend to visible regions
- Siamese CNN to get local feature maps in BEVs (point cloud + OSM)
- Adaptive Radial Fusion (ARF) for dynamic radial feature weighting
- Achieving +16% higher recall (KITTI, KITTI-360)













Outlook

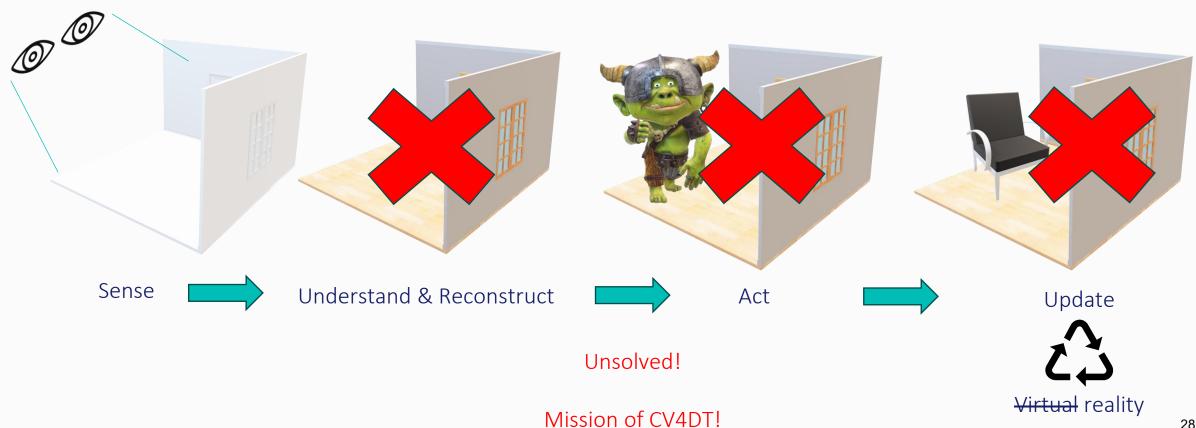








Digital Twins – they need 3D vision











3D computer vision for built environment

Can we go beyond the toy examples?





a plush dragon toy





A construction site







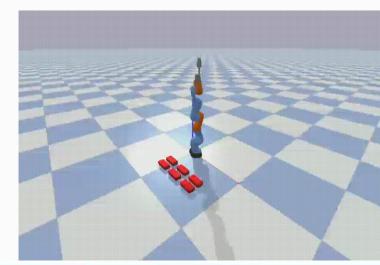


World Models – they need 3D vision

3D object-oriented semantic environment as physics prior to 2D/3D perception



Demo game – Driving in 3D reconstructed urban environment of Munich, Germany





Demo game – Walking in 3D reconstructed urban environment of Ingolstadt, Germany (toggle on and off improved models)

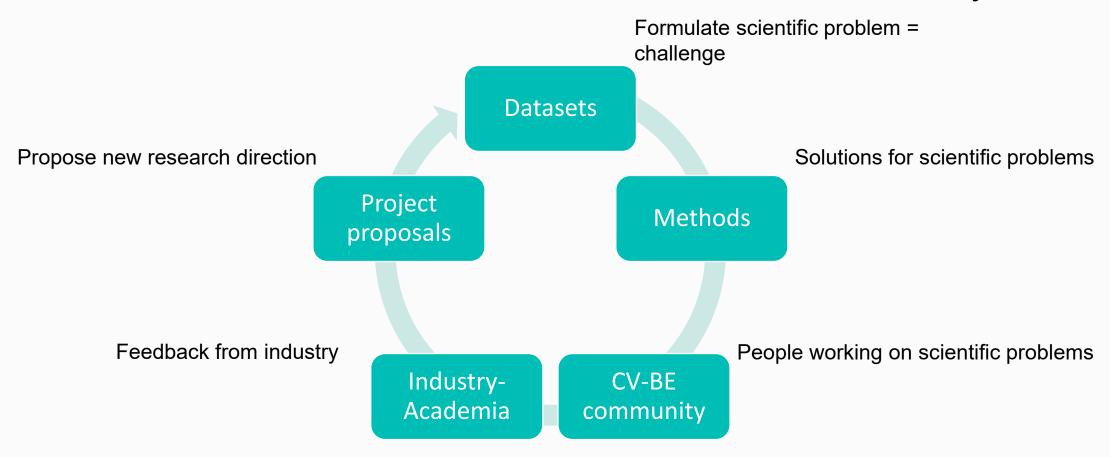








CV4DT – On the mission to create the built environment ecosystem











CamCV – Cambridge Community for Computer Vision

Gathering CV-related researchers of Cambridge!

- Webpage
- Regular meet-ups
- Invited talks
- Pre-conference events (e.g., CVPR/ICCV... paper presentations in Cambridge)
- ...and more?









CV4DT

