

[This is a static image]

# 16.485: VNAV - Visual Navigation for Autonomous Vehicles

**Luca Carlone**



Lecture 27: Research Directions in SLAM



<https://arxiv.org/abs/1606.05830>

## (Past,) Present, and Future of Simultaneous Localization And Mapping: Towards the Robust-Perception Age

Cesar Cadena, Luca Carlone, Henry Carrillo, Yasir Latif,  
Davide Scaramuzza, José Neira, Ian Reid, John J. Leonard

C. Cadena et al., "Past, Present, and Future of Simultaneous Localization and Mapping: Toward the Robust-Perception Age," in IEEE Transactions on Robotics, vol. 32, no. 6, pp. 1309-1332, Dec. 2016, doi: 10.1109/TRO.2016.2624754. © IEEE. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

**Abstract**—Simultaneous Localization And Mapping (SLAM) consists in the concurrent construction of a model of the environment (the *map*), and the estimation of the state of the robot moving within it. The SLAM community has made astonishing progress over the last 30 years, enabling large-scale real-world applications, and witnessing a steady transition of this technology to industry. We survey the current state of SLAM and consider future directions. We start by presenting what is now the *de-facto* standard formulation for SLAM. We then review related work, covering a broad set of topics including robustness and scalability in long-term mapping, metric and semantic representations for mapping, theoretical performance guarantees, active SLAM and exploration, and other new frontiers. This paper simultaneously serves as a position paper and tutorial to those who are users of SLAM. By looking at the published research with a critical eye, we delineate open challenges and new research issues, that still deserve careful scientific investigation. The paper also contains the authors' take on two questions that often animate discussions during robotics conferences: *Do robots need SLAM?* and *Is SLAM solved?*

**Index Terms**—Robots, SLAM, Localization, Mapping, Factor

### I. INTRODUCTION

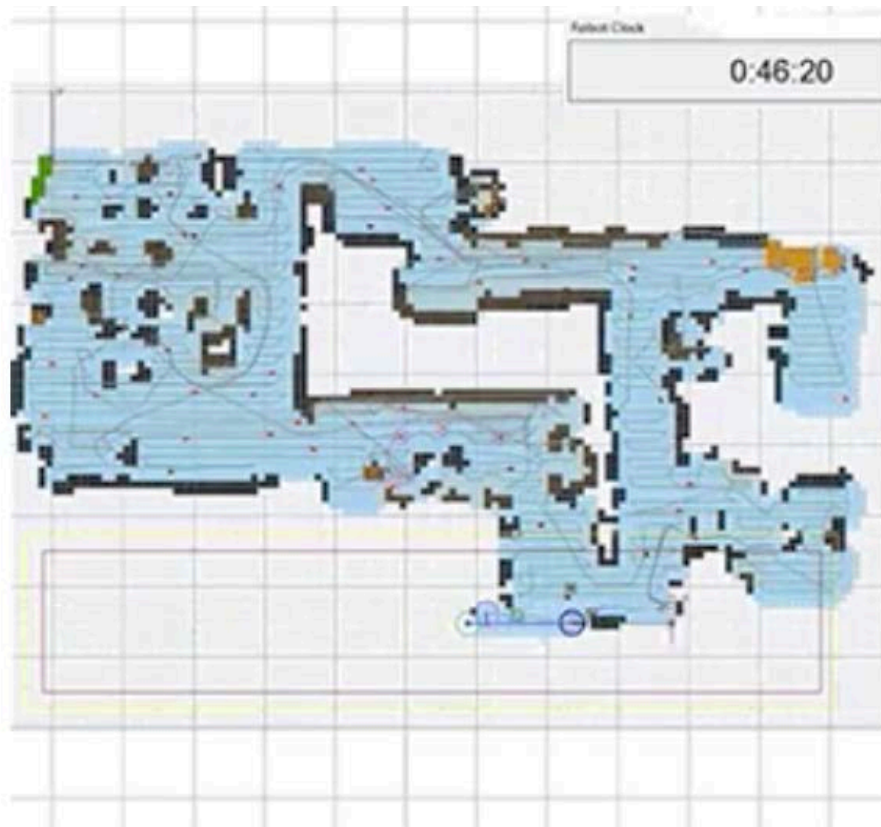
SLAM comprises the simultaneous estimation of the state of a robot equipped with on-board sensors, and the construction of a model (the *map*) of the environment that the sensors are perceiving. In simple instances, the robot state is described by its pose (position and orientation), although other quantities may be included in the state, such as robot velocity, sensor biases, and calibration parameters. The map, on the other hand, is a representation of aspects of interest (e.g., position of landmarks, obstacles) describing the environment in which the robot operates.

The need to use a map of the environment is twofold. First, the map is often required to support other tasks; for instance, a map can inform path planning or provide an intuitive visualization for a human operator. Second, the map allows limiting the error committed in estimating the state of the robot. In the absence of a map, dead-reckoning would quickly drift over time; on the other hand, using a map, e.g.,



# SLAM & SfM: Engineered Solutions / Applications

## Roomba 980 Vacuum Cleaner



## Kuka's Navigation Solution



## Mars Rovers (VO)



Source: public domain.



# SLAM & SfM: Engineered Solutions / Applications

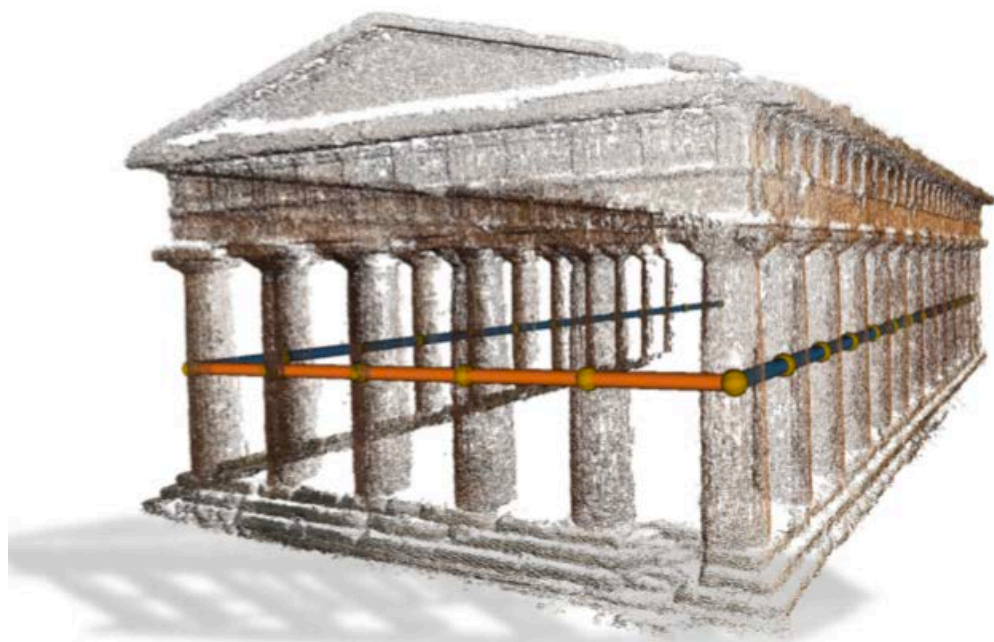
Precision agriculture



Google street view



Monitoring of historical sites



## Selling high-end property with drone mapping

USE CASES REAL ESTATE PIX4DMAPPER 3D MODELING MAPPING

22 OCTOBER 2015

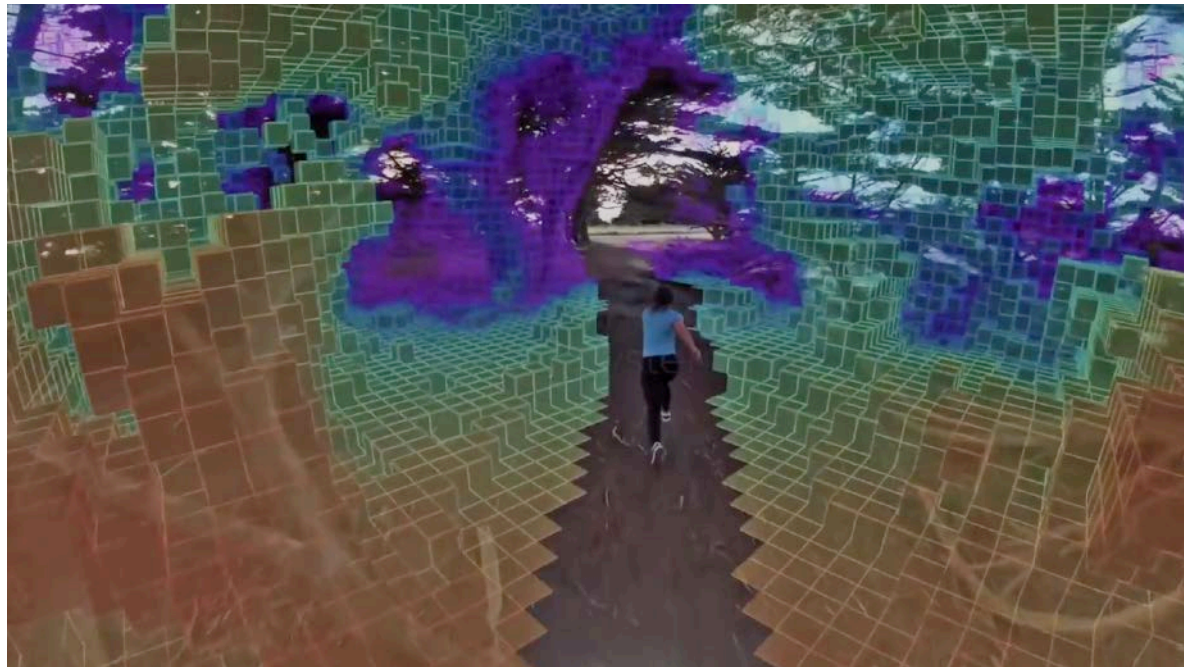
Pix4Dmapper is used to create a comprehensive 3D reconstruction of a luxury house for potential real estate clients.



# SLAM & SfM: Engineered Solutions / Applications

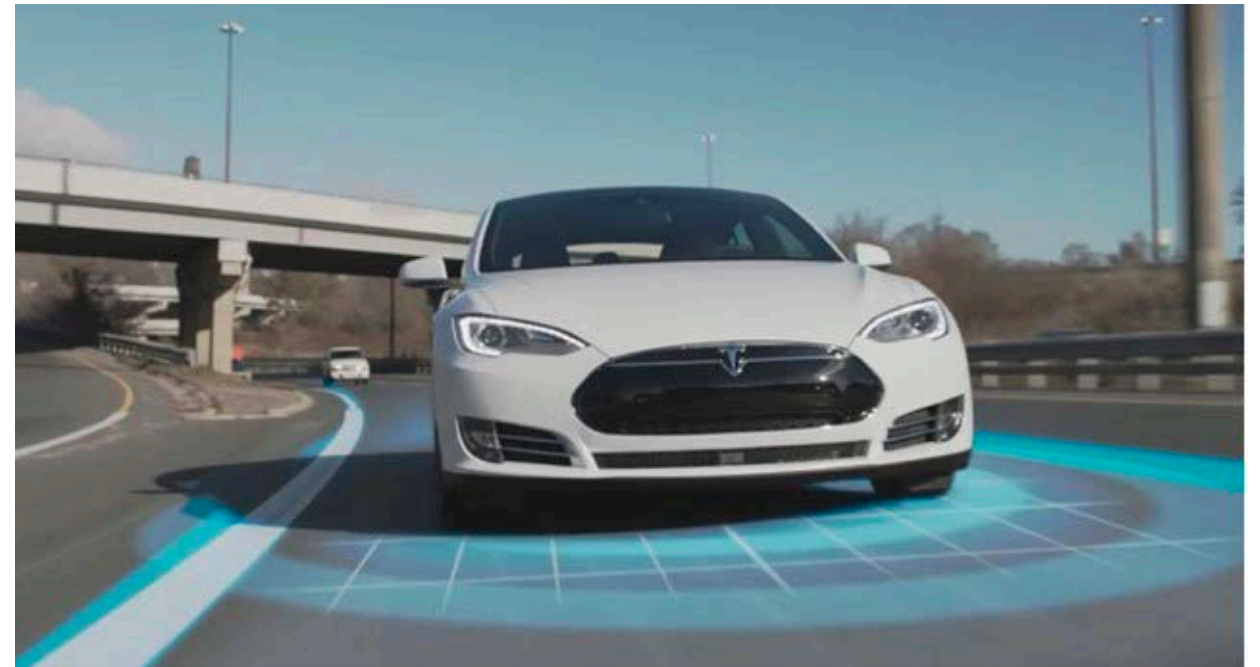
© Skydio, Inc. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

## Skydio R1 drone



© Tesla. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

## Tesla's autopilot



© Facebook. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

## Oculus Rift Goggles



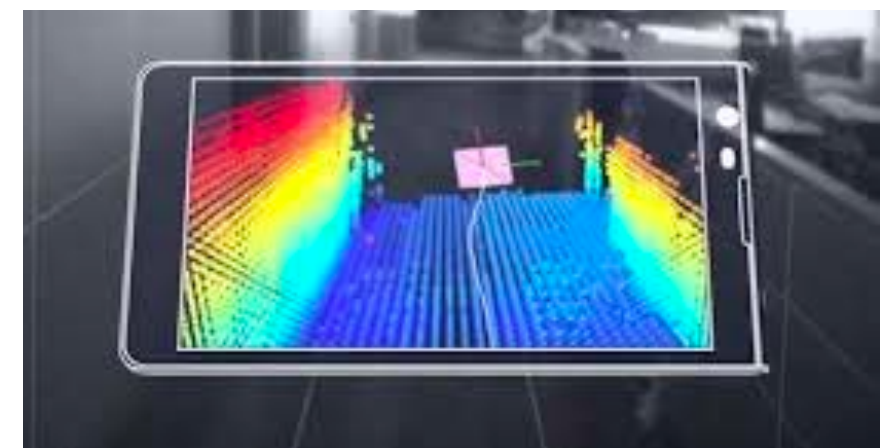
© Nintendo. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

## Pokemon Go



© Google. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

## Google Tango

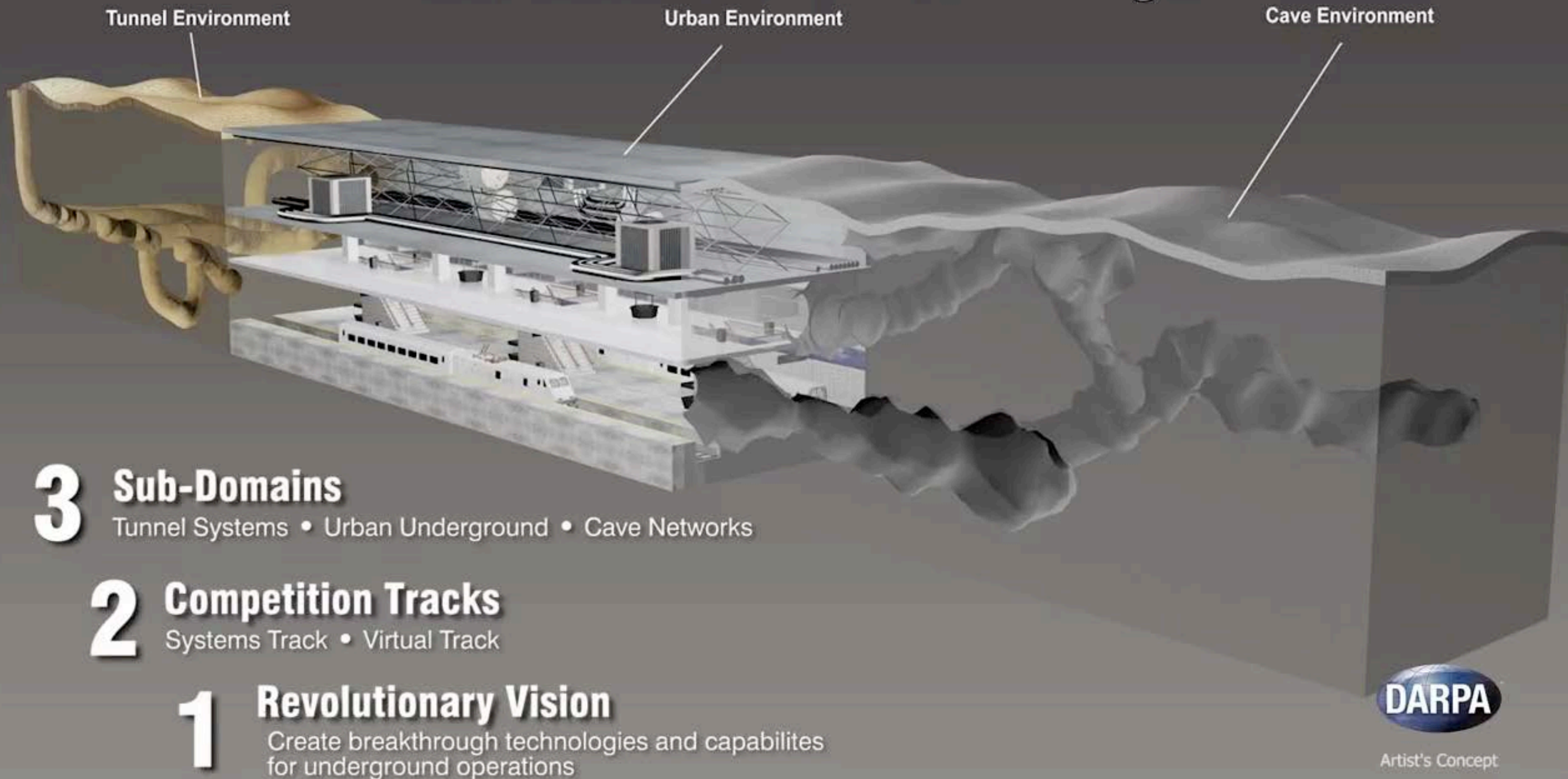


Reinvented as  
ARCore in 2017<sub>6</sub>



# SLAM & SfM: Engineered Solutions / Applications

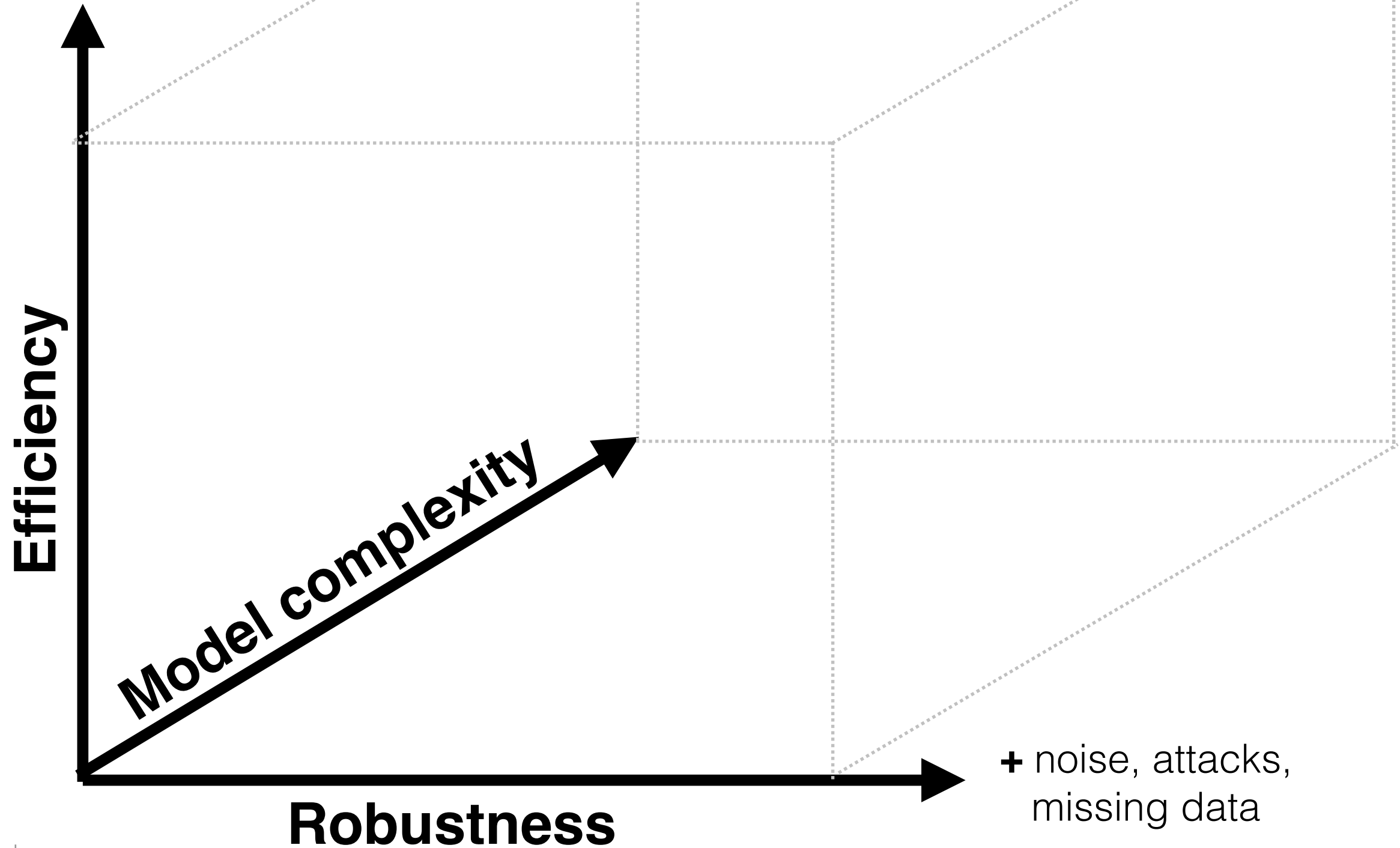
## DARPA Subterranean Challenge



DARPA Subterranean Challenge

# Axes of complexity

- power, size,  
time constants



# Axes of complexity

- power, size,  
time constants

**Efficiency**

**Model complexity**

**Robustness**

+ noise, attacks,  
missing data

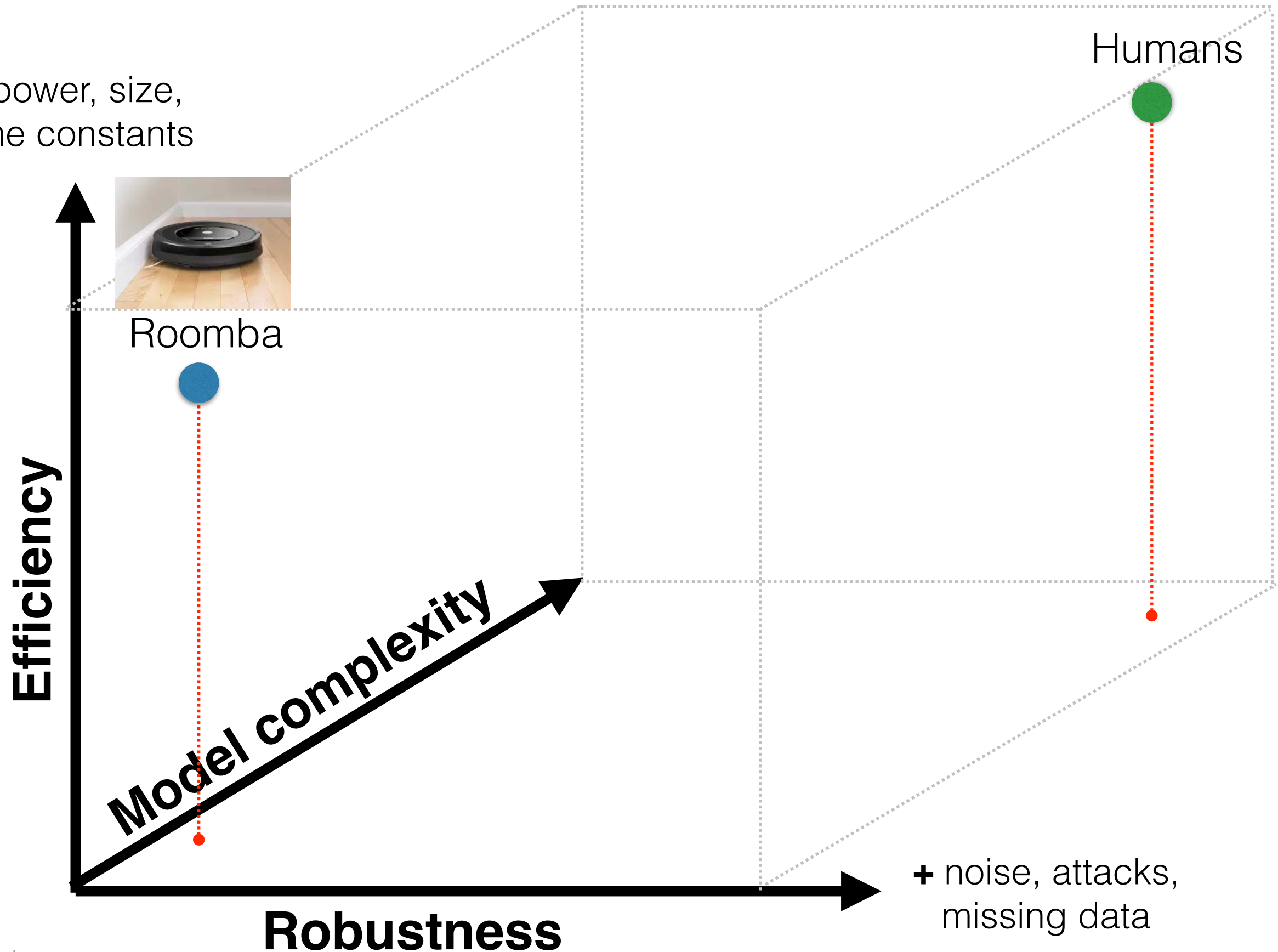
Humans





# Axes of complexity

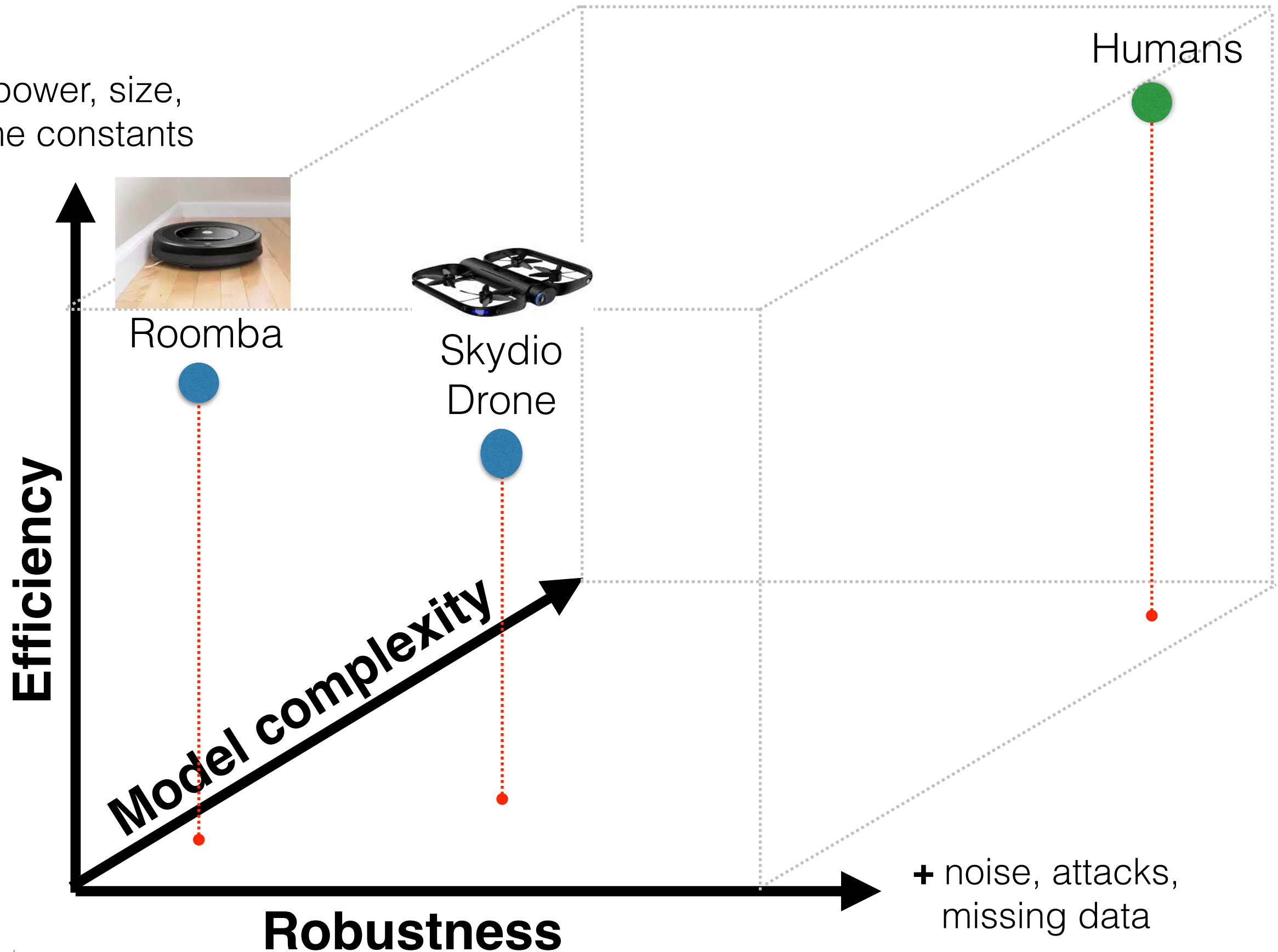
- power, size,  
time constants





# Axes of complexity

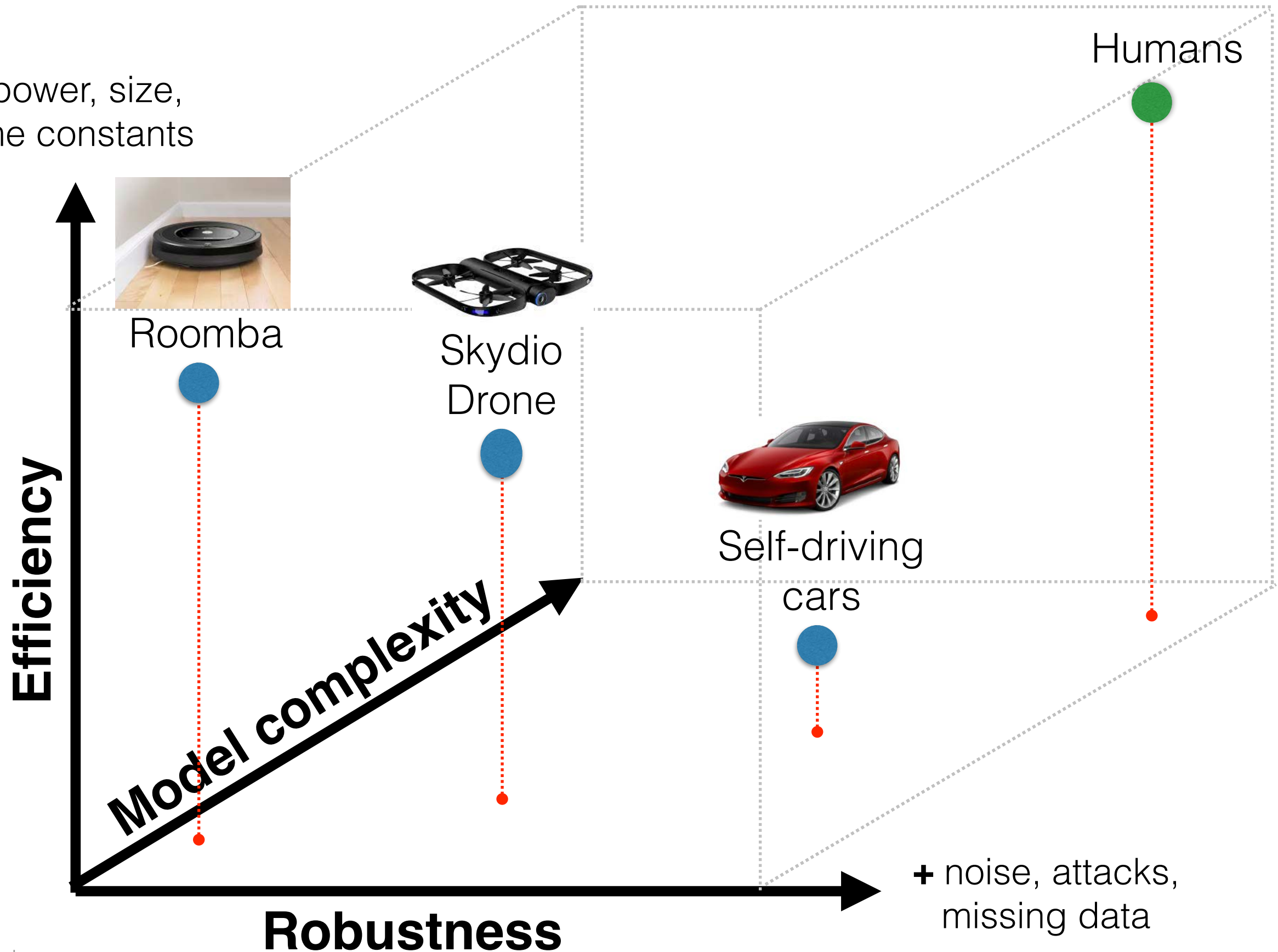
- power, size,  
time constants





# Axes of complexity

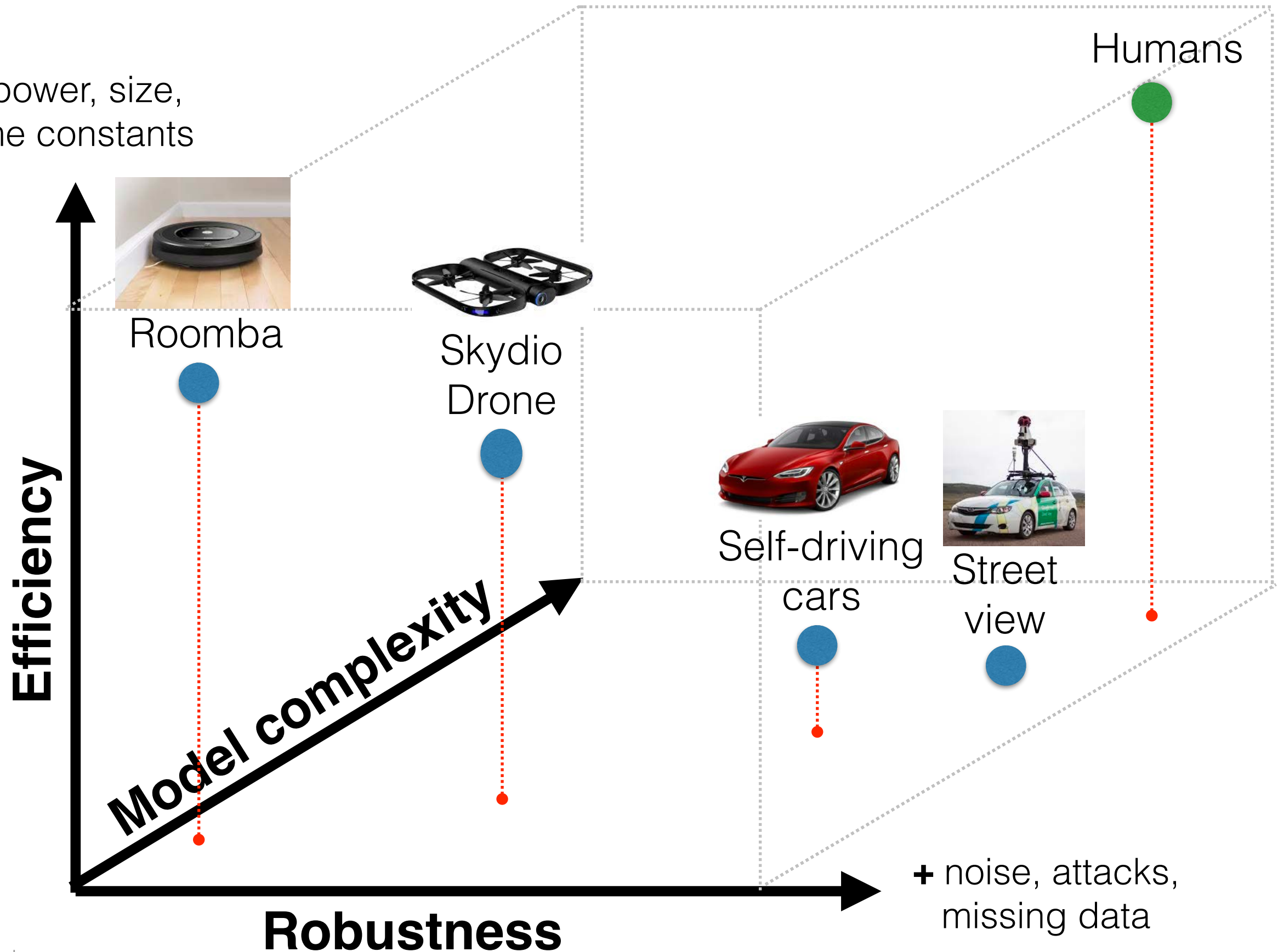
- power, size,  
time constants





# Axes of complexity

- power, size,  
time constants

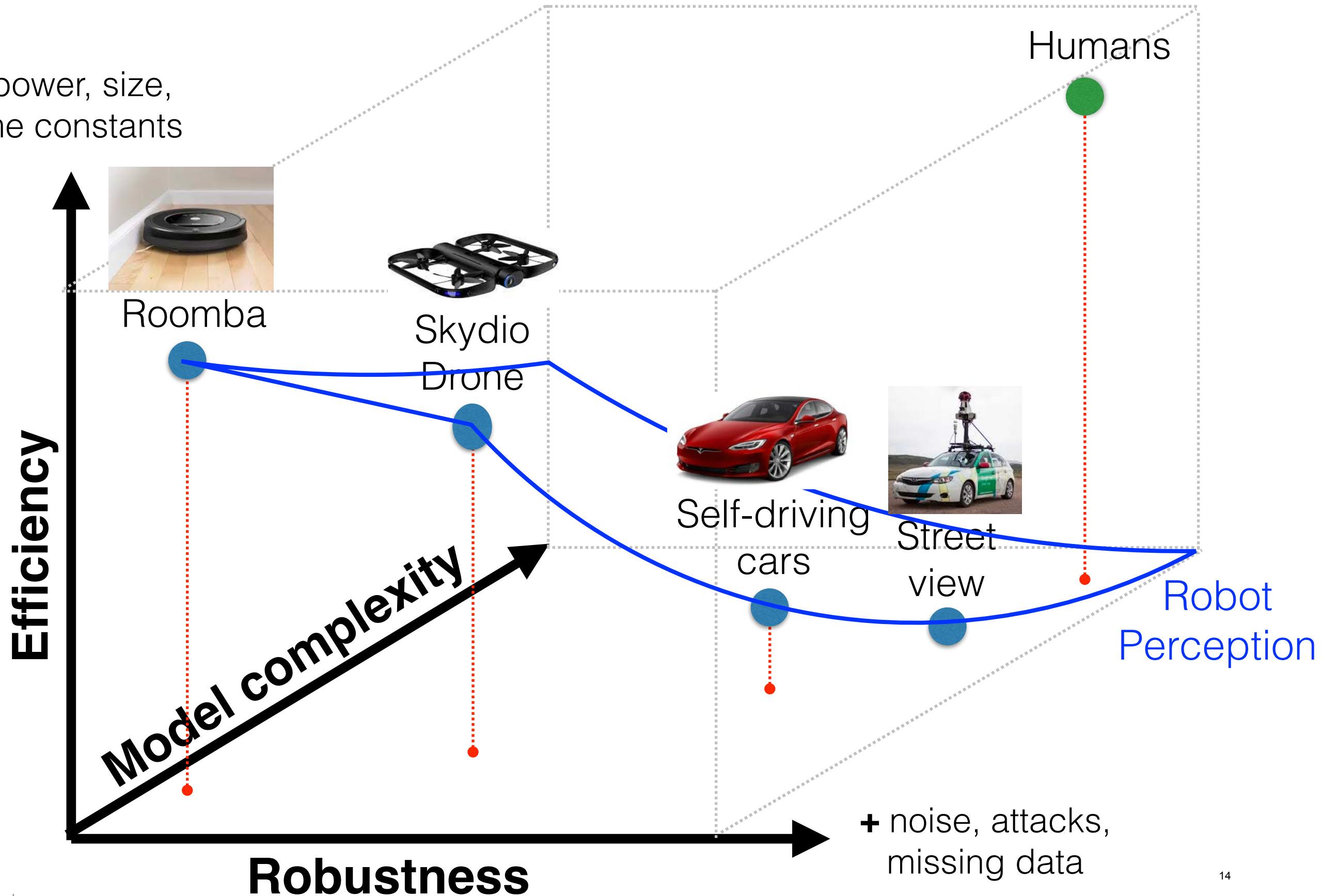


+ noise, attacks,  
missing data



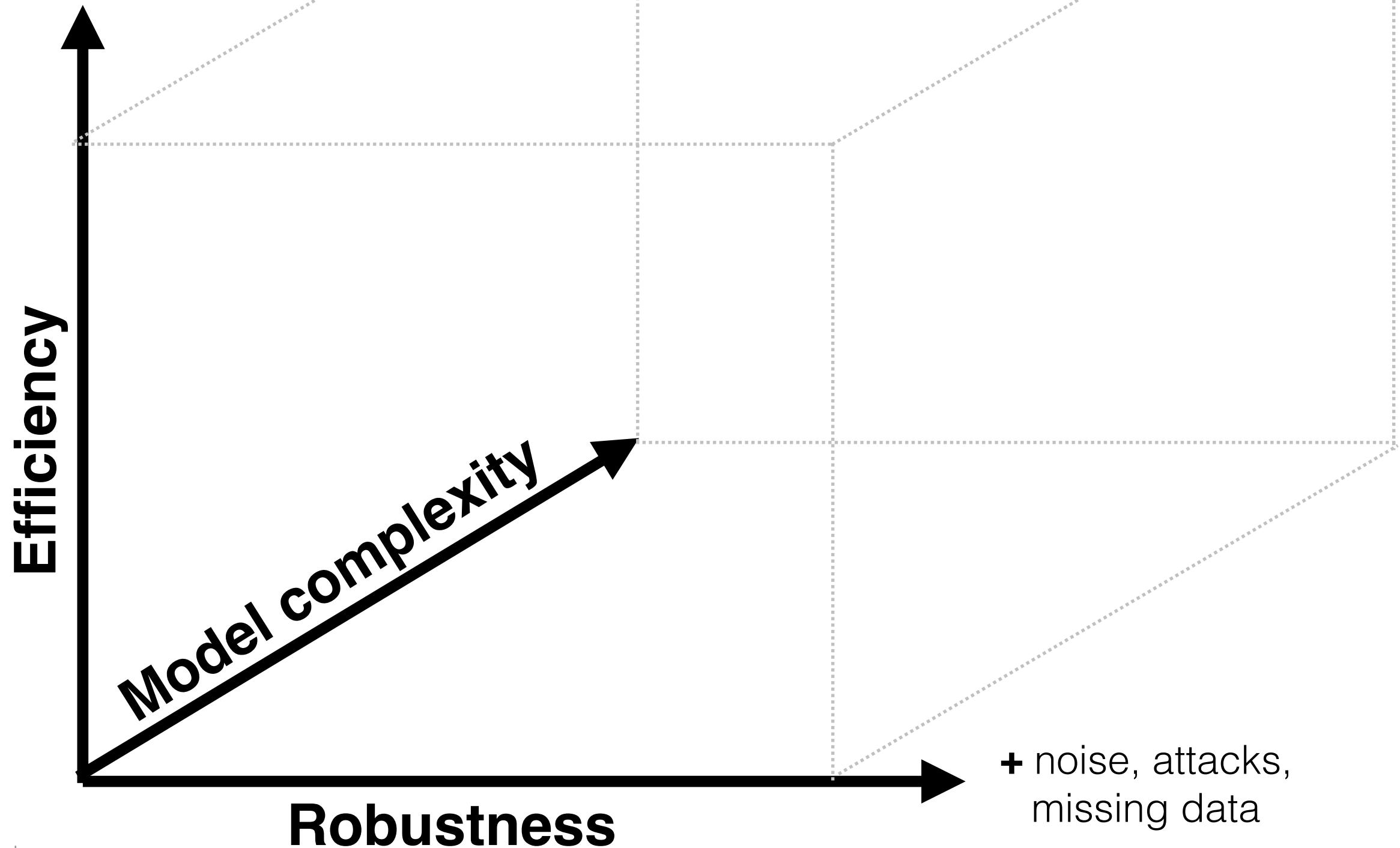
# Axes of complexity

- power, size,  
time constants



# Active Research Directions

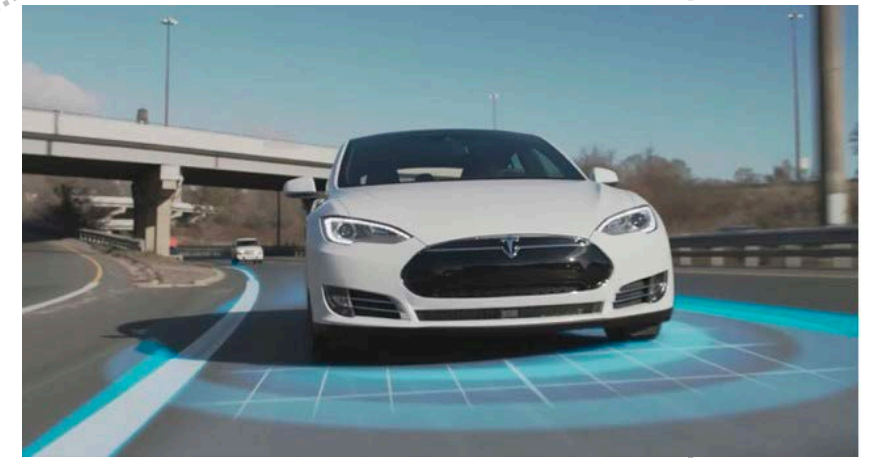
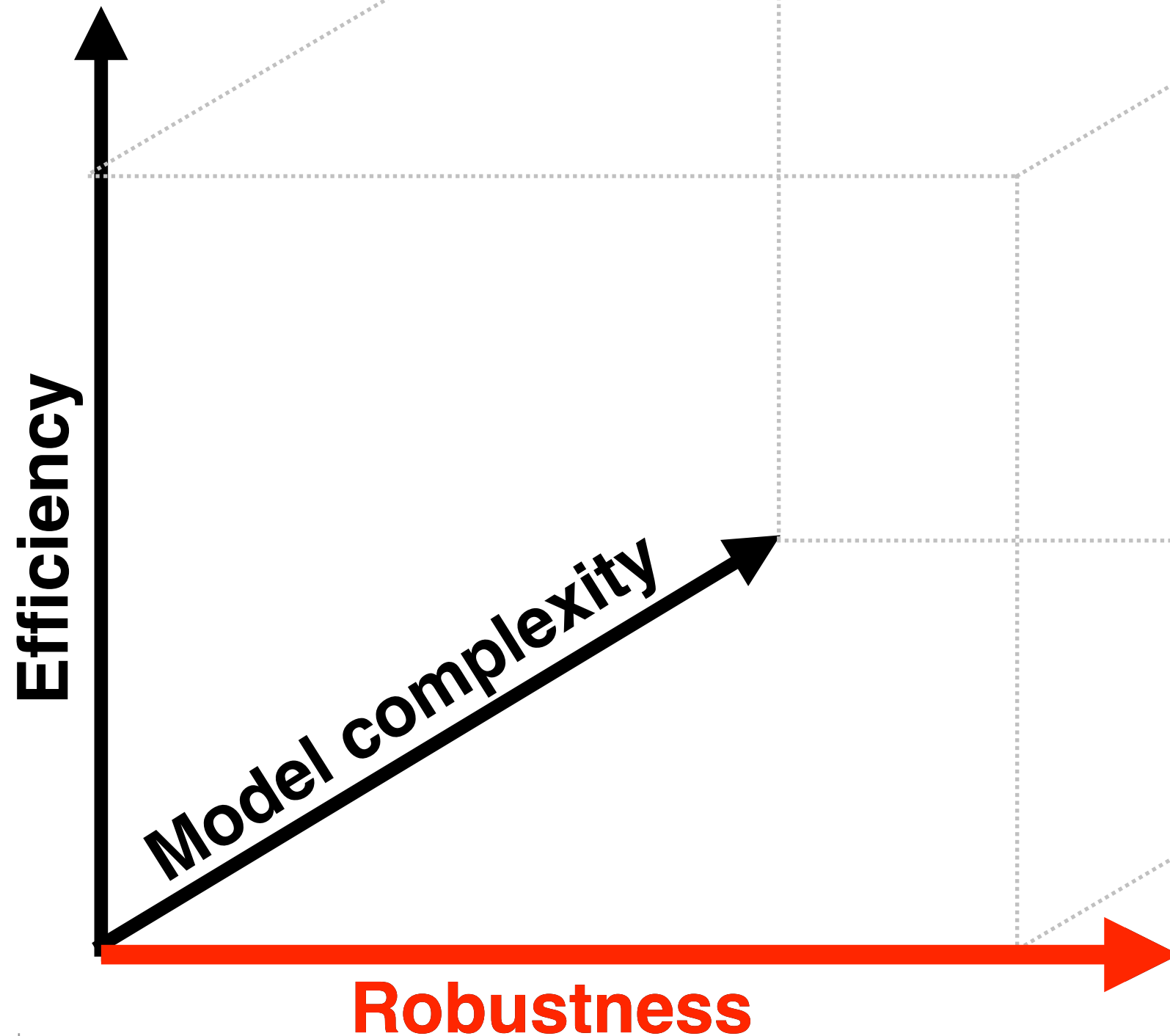
- power, size,  
time constants





# Active Research Directions

- power, size,  
time constants

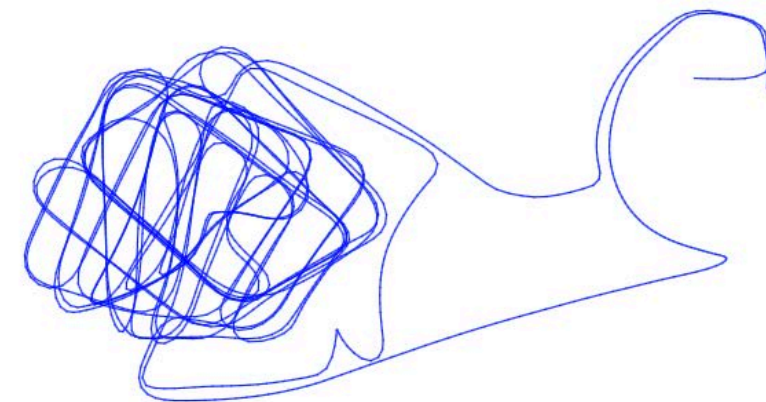
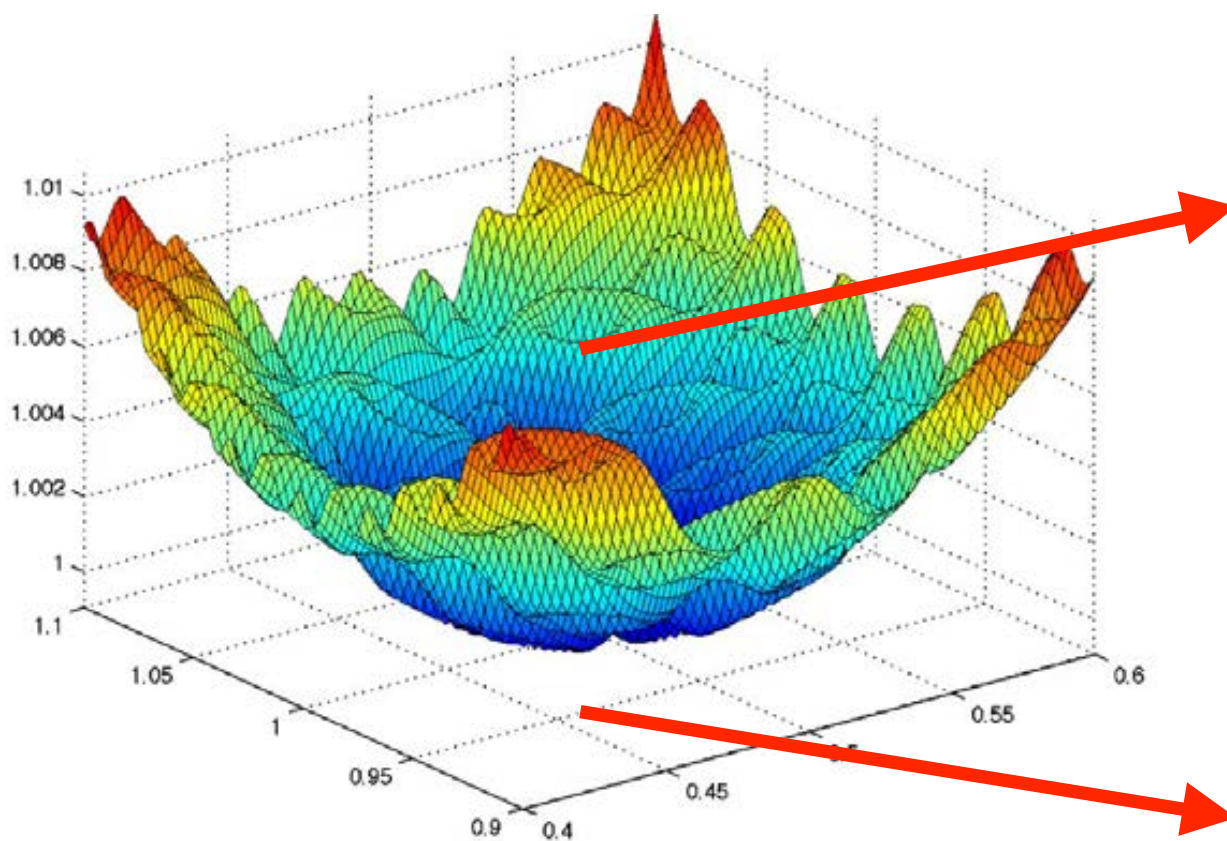


+ noise, attacks,  
missing data

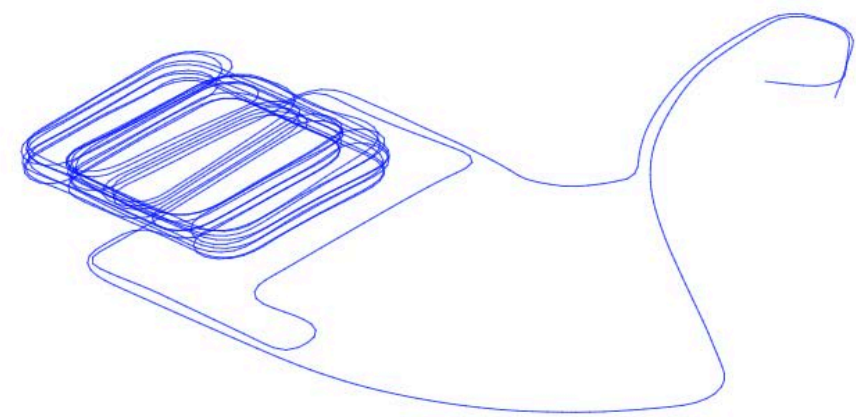
# Robustness to noise

$$\min_{\substack{\{p_i \in \mathbb{R}^3\} \\ \{R_i \in \text{SO}(3)\}}} \sum_{(i,j) \in \mathcal{E}} \frac{1}{\sigma_p^2} \left\| \bar{p}_{ij} - R_i^\top (p_j - p_i) \right\|^2 + \frac{1}{\sigma_R^2} \left\| \bar{R}_{ij} - R_i^\top R_j \right\|_F^2$$

**non-convex optimization**



Suboptimal critical point



Optimal estimate

**Iterative methods (e.g., gradient descent)  
may get stuck into bad minima**

**Initial guess gets worse when noise is large**

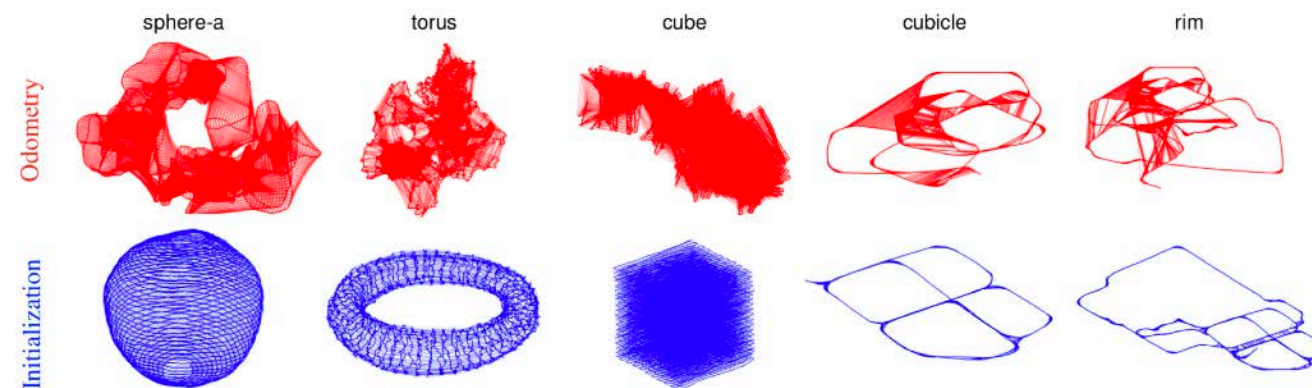


# Robustness to **noise** (convergence)

- **Analysis:** number of minima, basin of attraction of iterative solvers (Gauss-Newton), factors impacting quality of solution

Initialization Techniques for 3D SLAM: a Survey on Rotation Estimation and its Use in Pose Graph Optimization

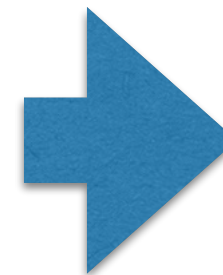
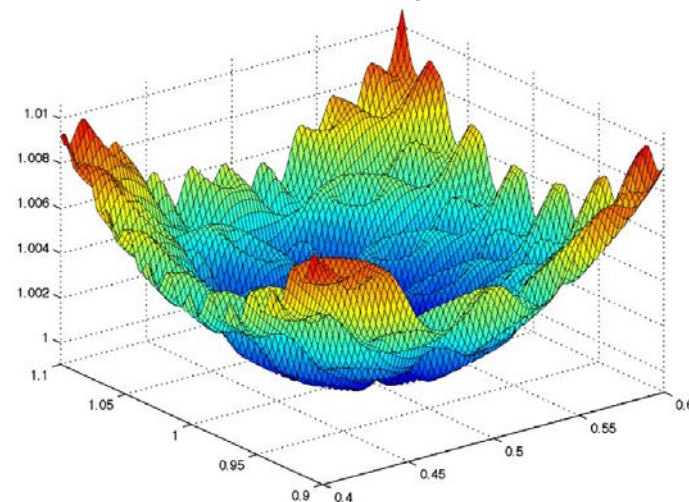
Luca Carlone, Roberto Tron, Kostas Daniilidis, and Frank Dellaert



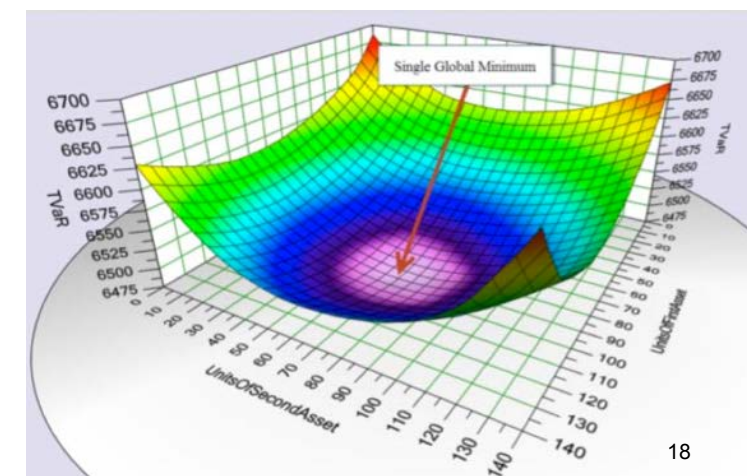
- **Initialization Techniques**

- **Global Solvers**

non convex problem



convex relaxation

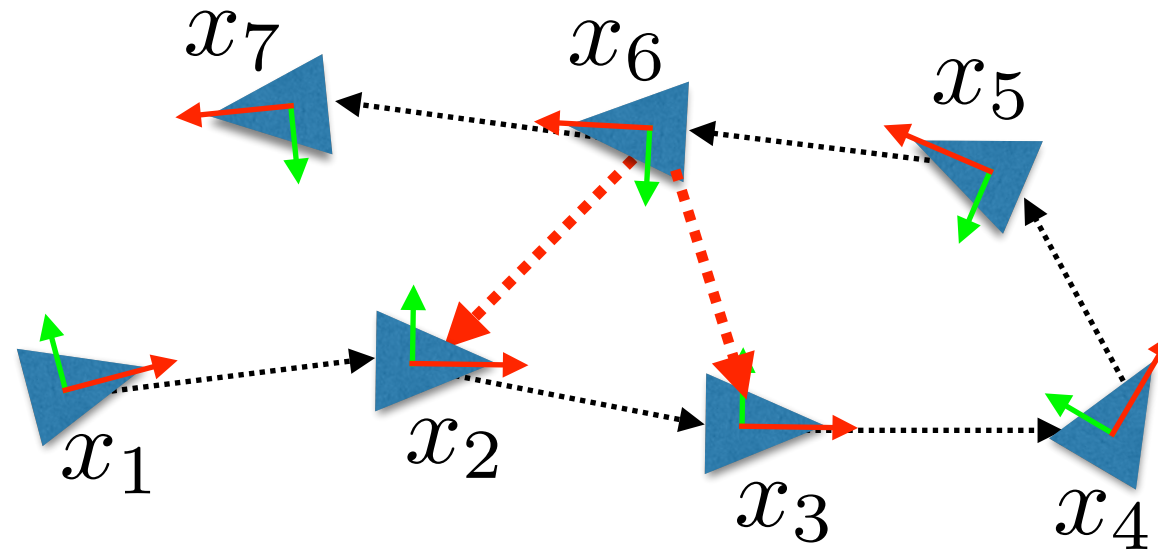


# What if place recognition fails?

---

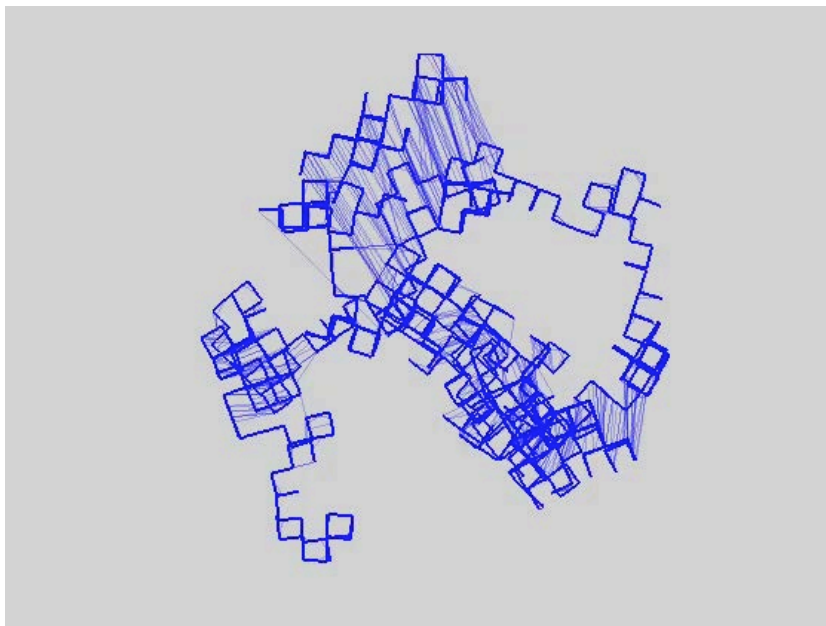


# What if place recognition fails?

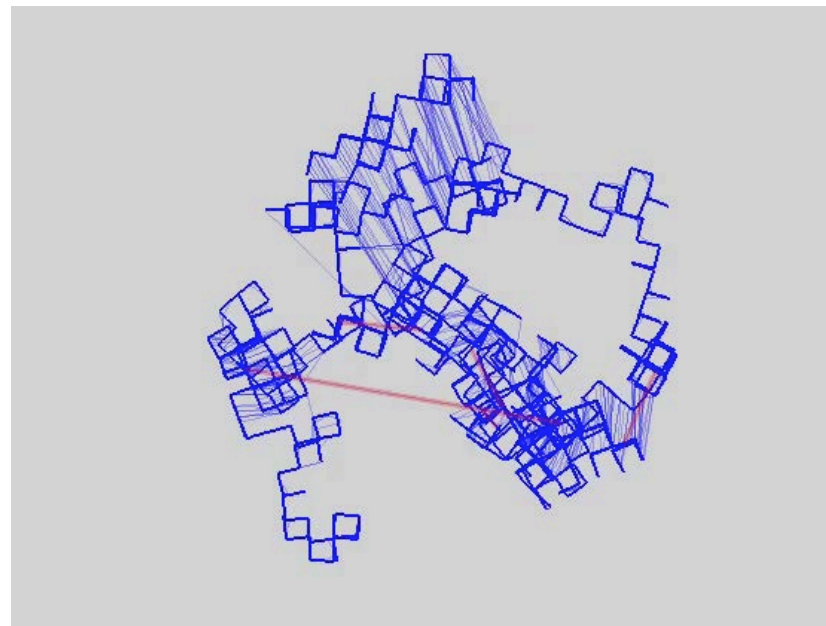


**outliers:** completely incorrect measurements  
(Perceptual Aliasing)

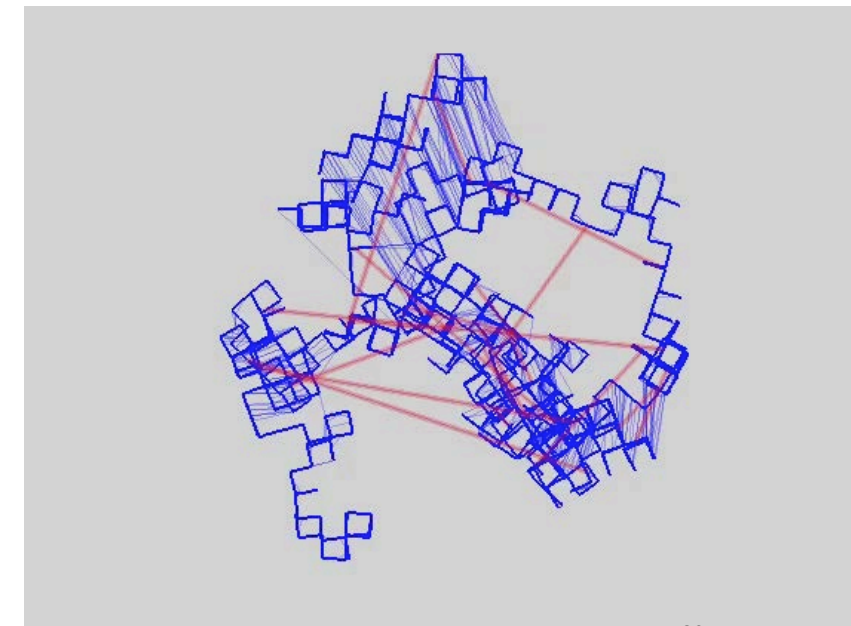
0 outliers



5 outliers



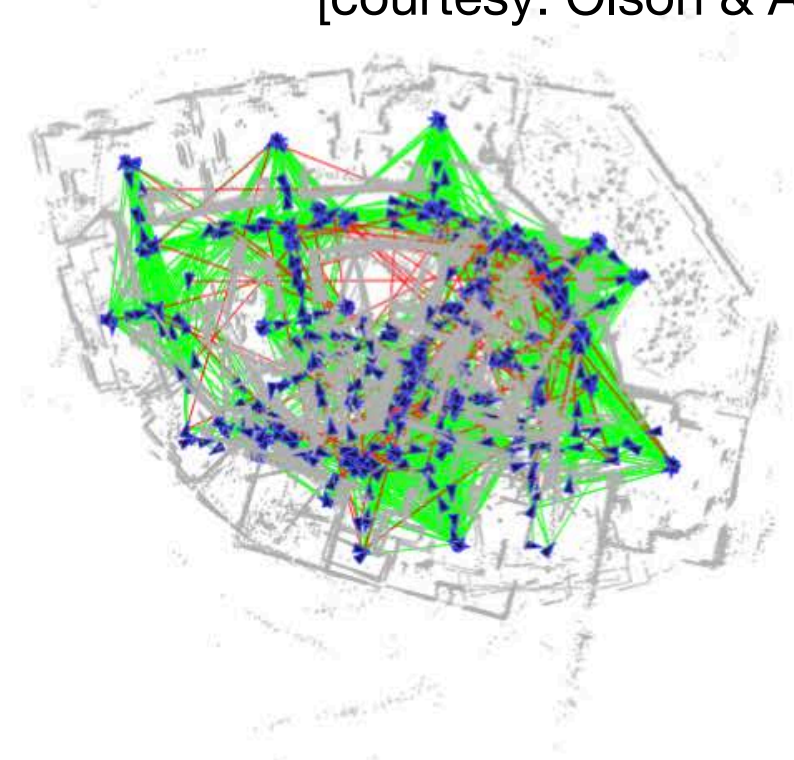
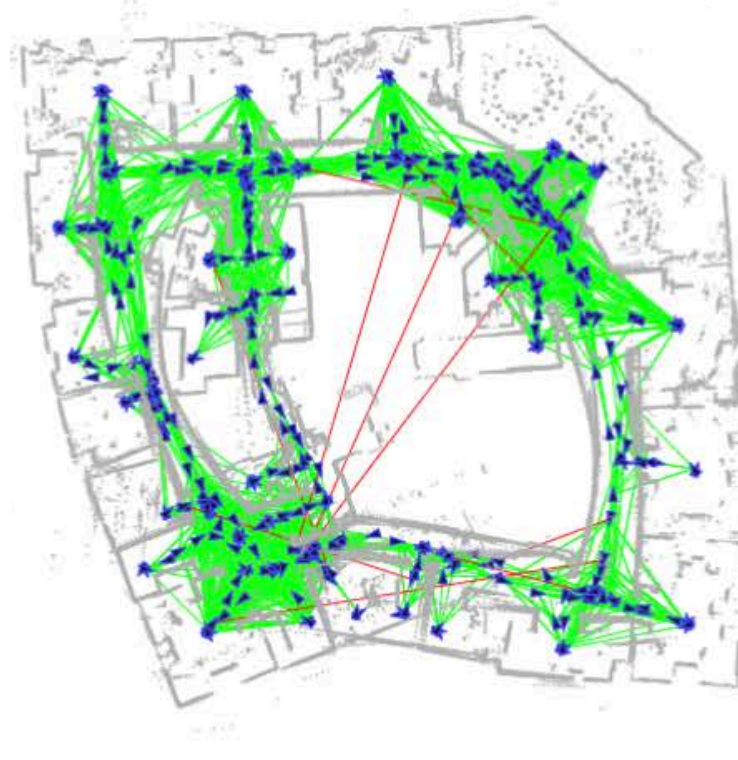
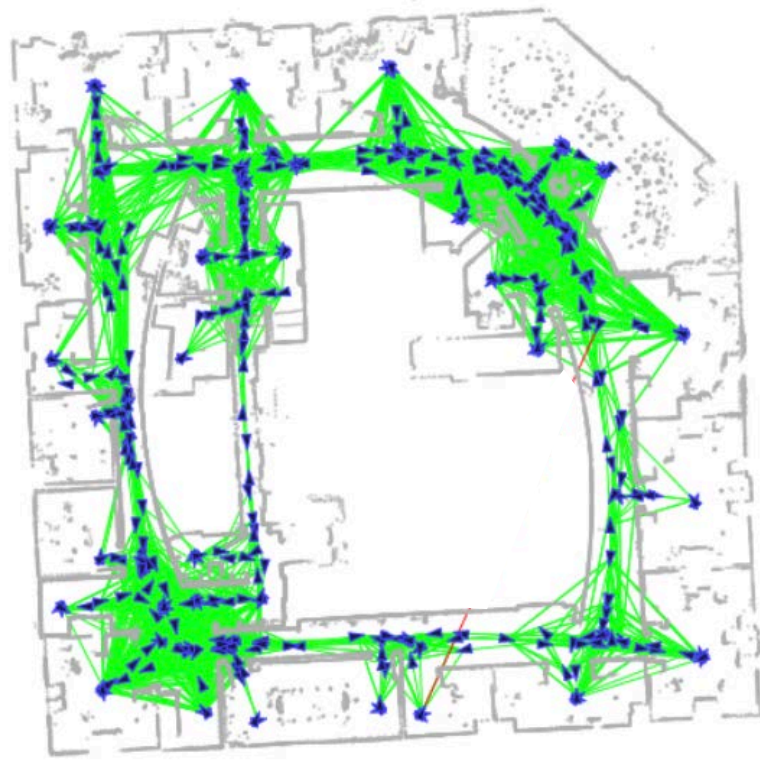
20 outliers





# Robustness to outliers

[courtesy: Olson & Agarwal]



correct but noisy measurements



outliers (completely wrong measurements)

© Olson and Agarwal. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

least-square SLAM estimators catastrophically fail if outliers are not carefully handled



“Google employs a small army of human operators to manually check and correct the maps” [Wired]



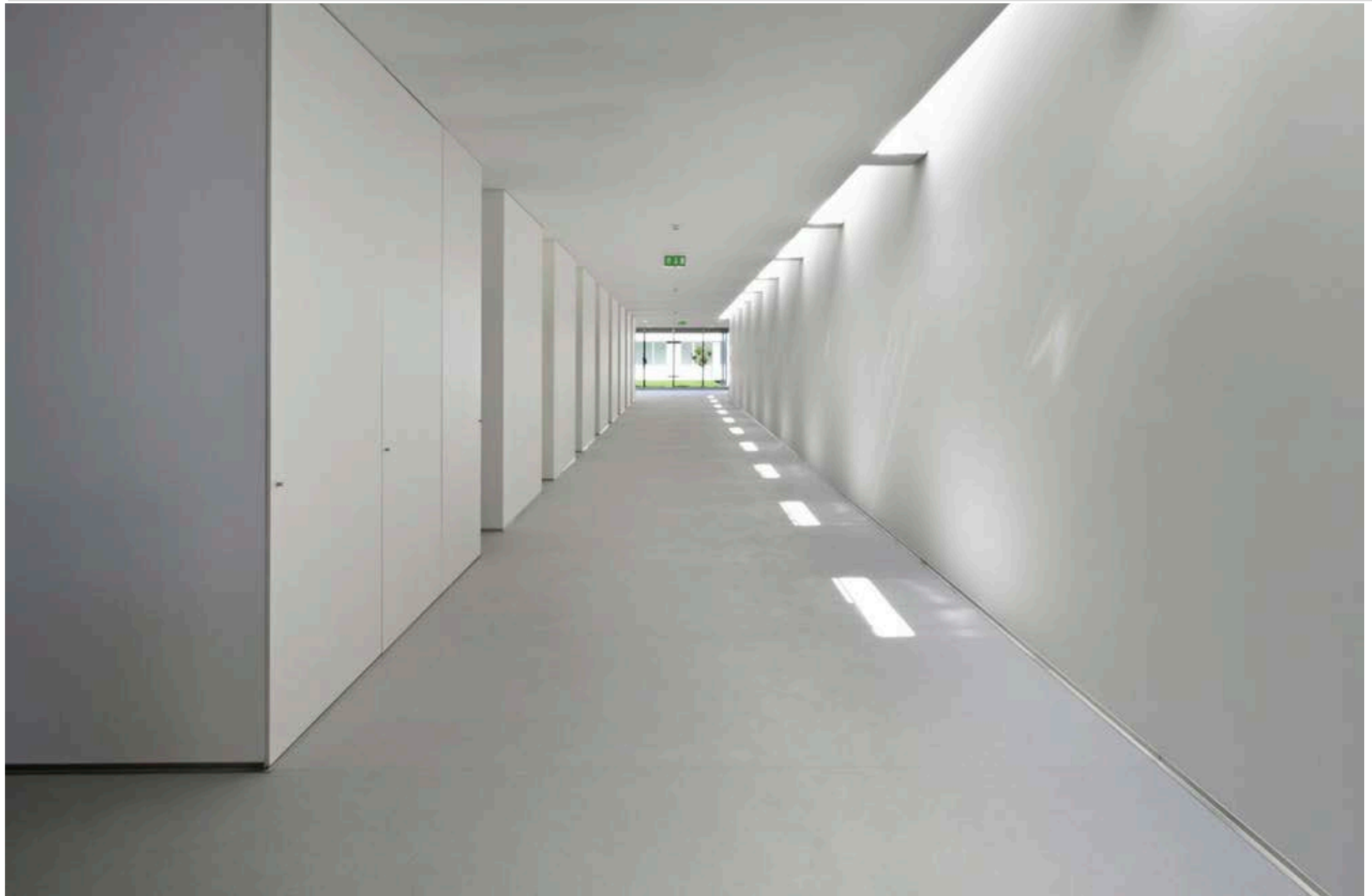
# Robustness to **dynamic scenes**

---





# Robustness to **missing data**



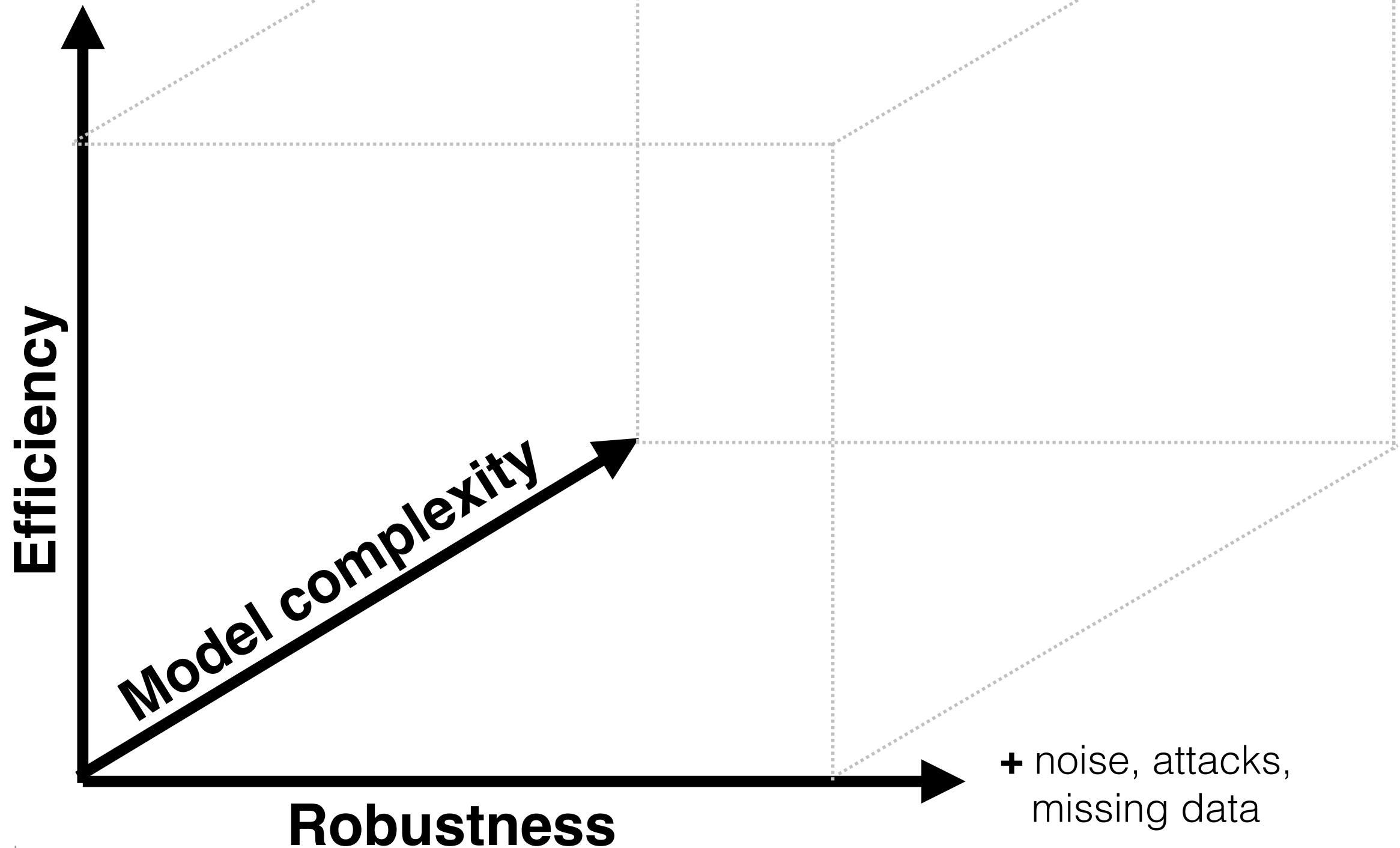
. Zhang, M. Kaess and S. Singh, "On degeneracy of optimization-based state estimation problems," 2016 IEEE International Conference on Robotics and Automation (ICRA), Stockholm, Sweden, 2016, pp. 809-816, doi: 10.1109/ICRA.2016.7487211 © IEEE All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

[Zhang, Kaess, Singh, On Degeneracy of Optimization-based State Estimation]



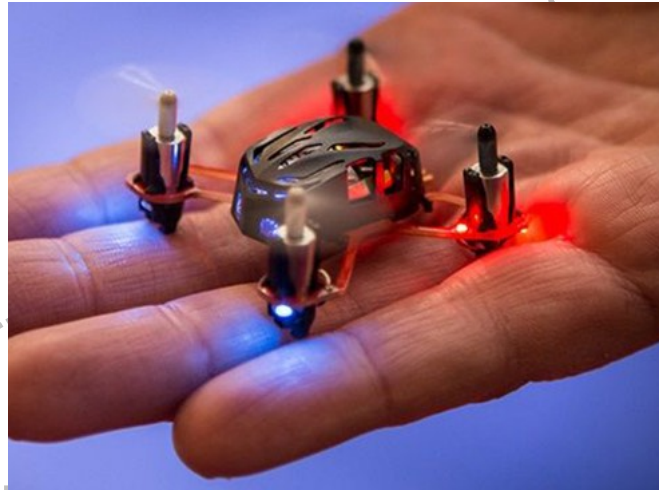
# Active Research Directions

- power, size,  
time constants



# Active Research Directions

- power, size,  
time constants



**Efficiency**

**Model complexity**

**Robustness**

+ noise, attacks,  
missing data

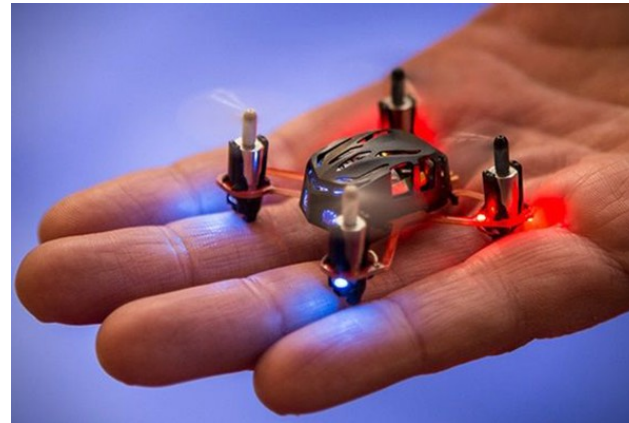


# Efficiency and Miniaturization

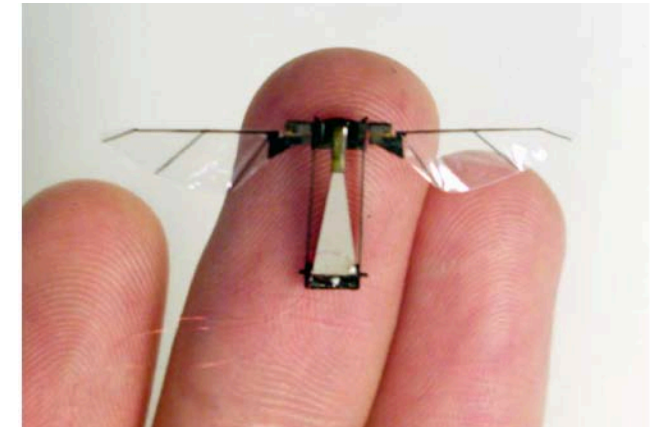
> 10 W



< 5 W



< 200 mW



## Human vision



## Machine vision



**Data stream**

$10^8 - 10^9$  bits/second

$5 \cdot 10^8$  bits/second (stereo)

**Performance**

parse scene: 13ms

object detection: 22ms (GPU)

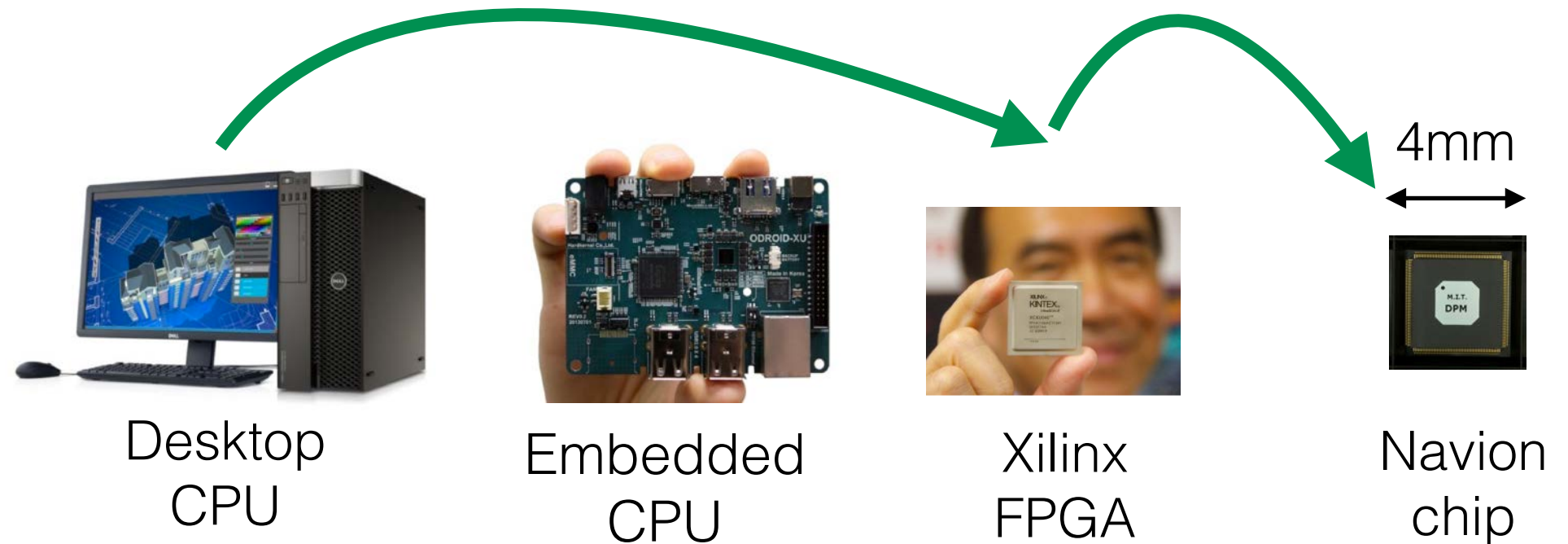
SLAM: >100ms

**Power**

20W

250 W (Titan X GPU)<sub>26</sub>

# Algorithms-and-hardware co-design



latency  
image processing

50 ms

200ms

50 ms

**22ms**

latency  
MAP estimation

80 ms

400ms

200ms

**30ms**

power

26.1 W

2.33 W

1.46 W

**24mW**

accuracy

16cm

16cm

19cm

23cm



# Efficiency and Miniaturization

---





# Active Research Directions

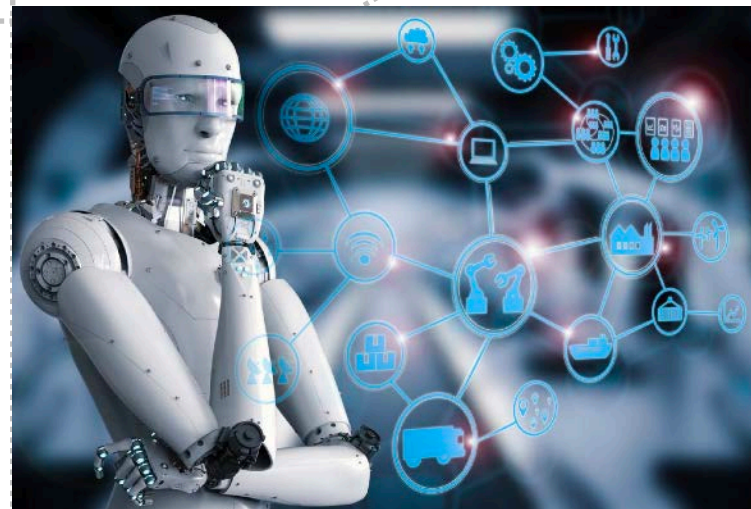
- power, size,  
time constants

Efficiency

Model complexity

Robustness

+ noise, attacks,  
missing data

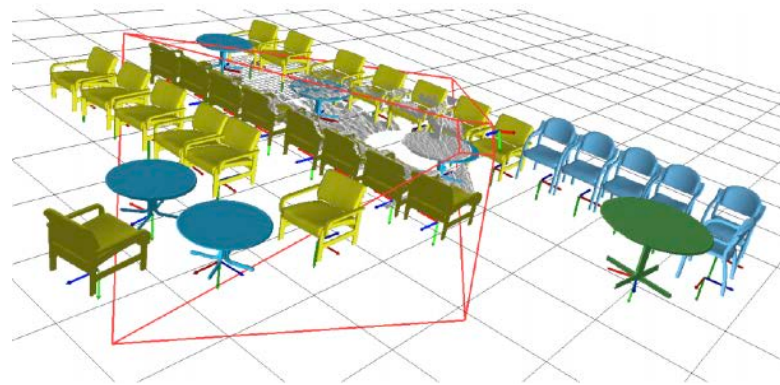
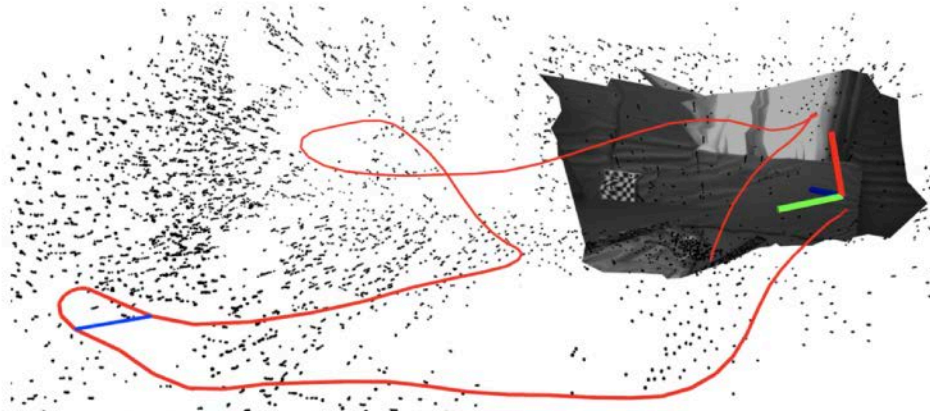




# Mind the Gap with Human Perception



Sparse  
or  
Dense  
Point  
Clouds  
Object-  
based  
Maps  
[Salas-  
Moreno,  
CVPR'13]



.. lines, voxels, meshes

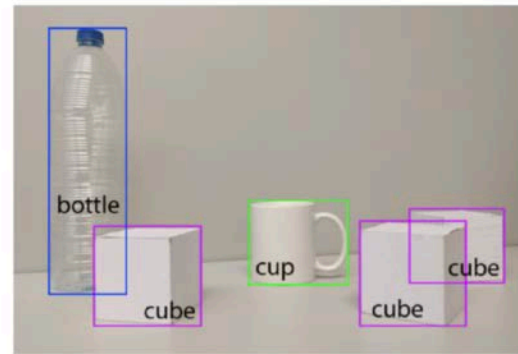




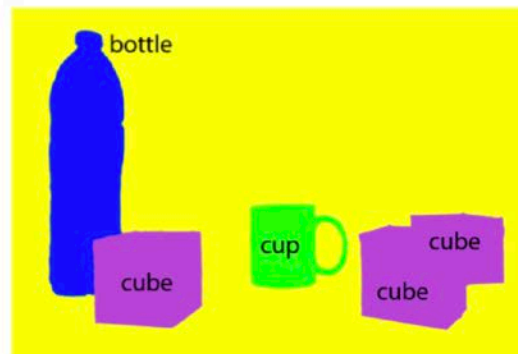
# High-level Understanding: Opportunities



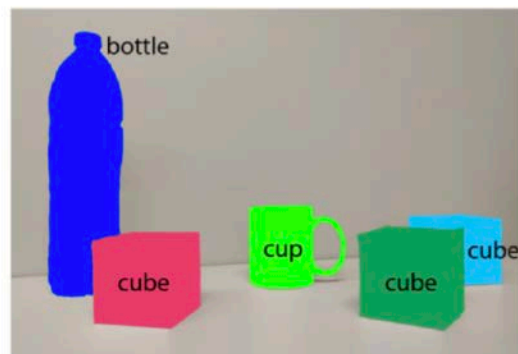
(a) Image classification



(b) Object localization

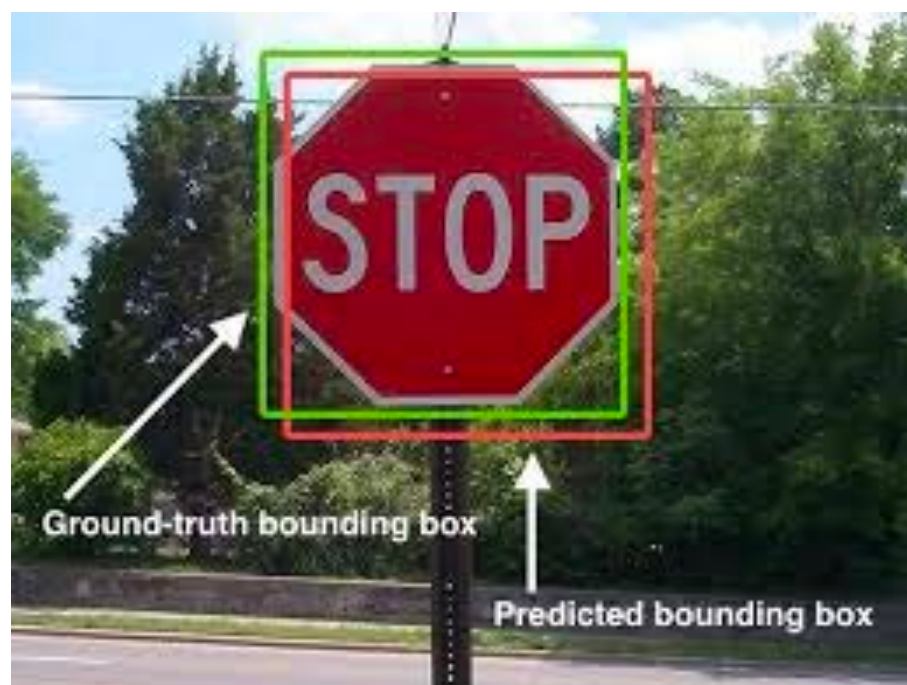


(c) Semantic segmentation



(d) Instance segmentation

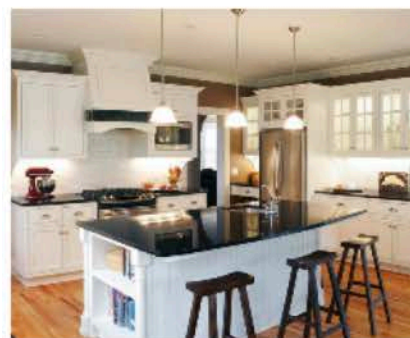
[Garcia-Garcia et al., 2017]



## 2.1. COCO Detection Challenge



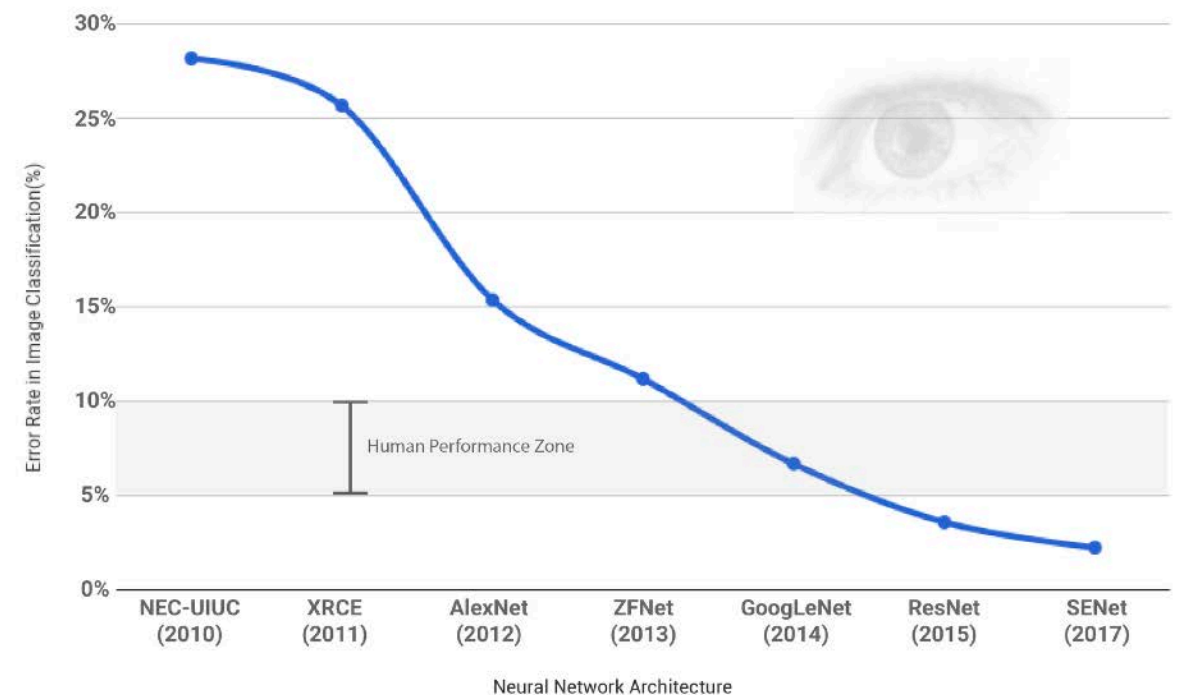
## 3. Places Challenges



Scene Parsing



Instance Segmentation

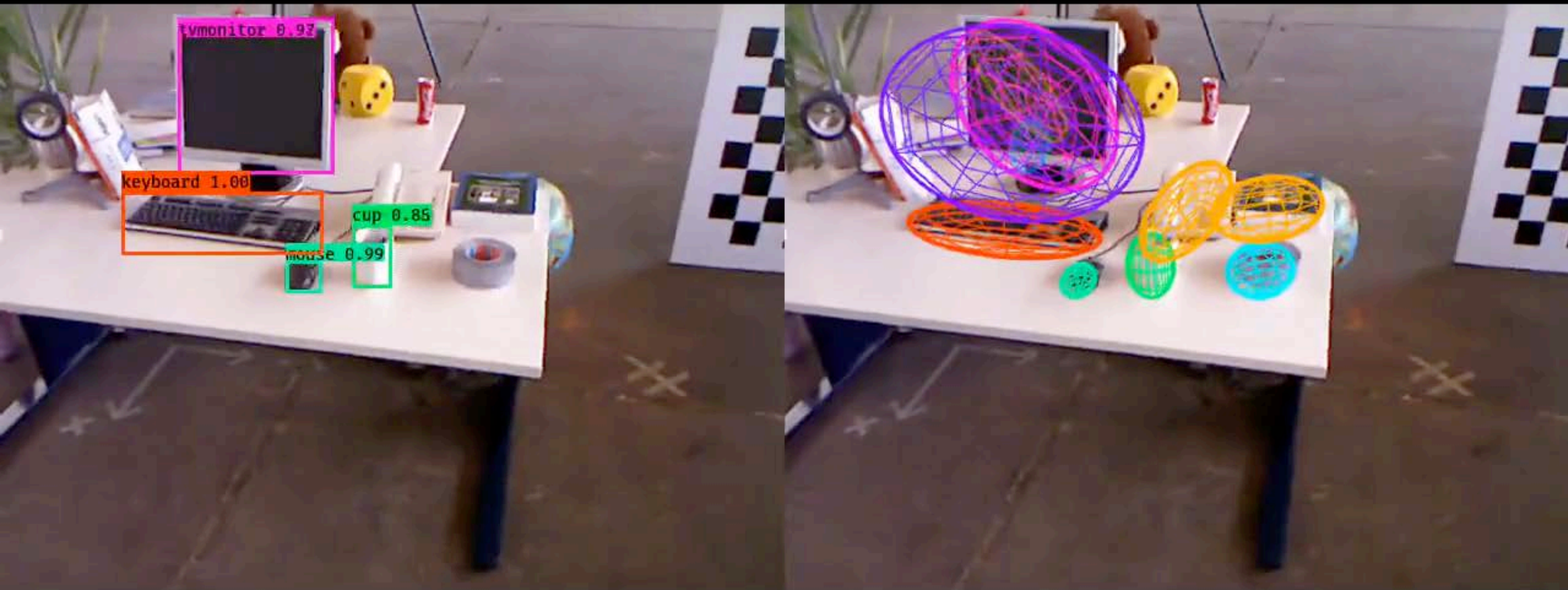


The deep learning revolution!



# Sparse Object-level SLAM

## QuadricSLAM



**Raw Detections**

**Projected Landmarks**

Figure 3 in Lachlan Nicholson, Michael Milford, and Niko Sunderhauf, "QuadricSLAM: Dual Quadrics from Object Detections as Landmarks in Object-oriented SLAM." IEEE ROBOTICS AND AUTOMATION LETTERS © IEEE. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

[Sunderhauf and Milford, 2017]

# Dense Metric-Semantic SLAM on GPU

---

## SemanticFusion: Dense 3D Semantic Mapping with Convolutional Neural Networks

John McCormac, Ankur Handa,  
Andrew Davison, Stefan Leutenegger

Dyson Robotics Lab, Imperial College London



# Dense Metric-Semantic SLAM on CPU

Kimera-Semantics



Kimera-RPGO



Kimera-VIO & Mesher



## Kimera



Top down view

Fast multi-threaded open-source code: <https://github.com/MIT-SPARK/Kimera>

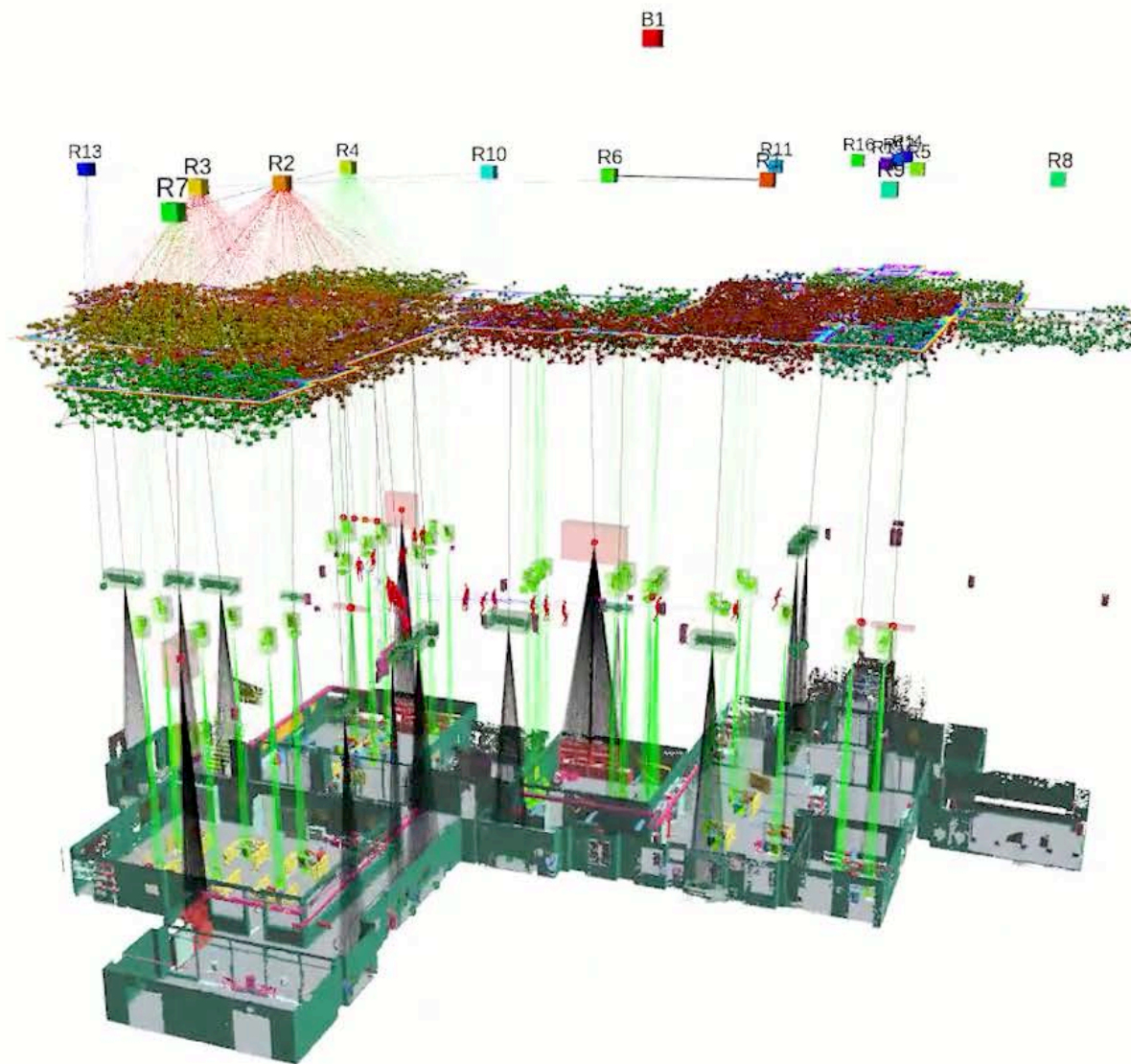


Speed

A. Rosinol, M. Abate, Y. Chang, L. Carlone, Kimera: an Open-Source Library for Real-Time Metric-Semantic Localization and Mapping. IEEE Intl. Conf. on Robotics and Automation (ICRA), 2020. arXiv:1910.02490 © IEEE. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

Rosinol, Abate, Chang, Carlone. Kimera: an open-source library for real-time metric-semantic localization and mapping. ICRA 2020.

# High-level Understanding: 3D Scene Graphs



- Directed graph, where:
- **nodes** are *spatial concepts* (i.e., concepts grounded in 3D)
  - **edges** represent spatio-temporal relations between concepts (e.g., agent “i” in room “j” at time t)

**We present a unified representation for actionable spatial perception:  
3D Dynamic Scene Graphs (DSGs)**

Figure 1 in Antoni Rosinol et al, "3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans." Robotics: Science and Systems 2020 Corvallis, Oregon, USA, July 12-16, 2020 © Antoni Rosinol et al. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

[Armeni et al., 3D scene graph: A structure for unified semantics, 3D space, and camera. ICCV'19]

[Rosinol et al., 3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans, RSS'20]



# High-level Understanding: 3D Scene Graphs

RGB Frame

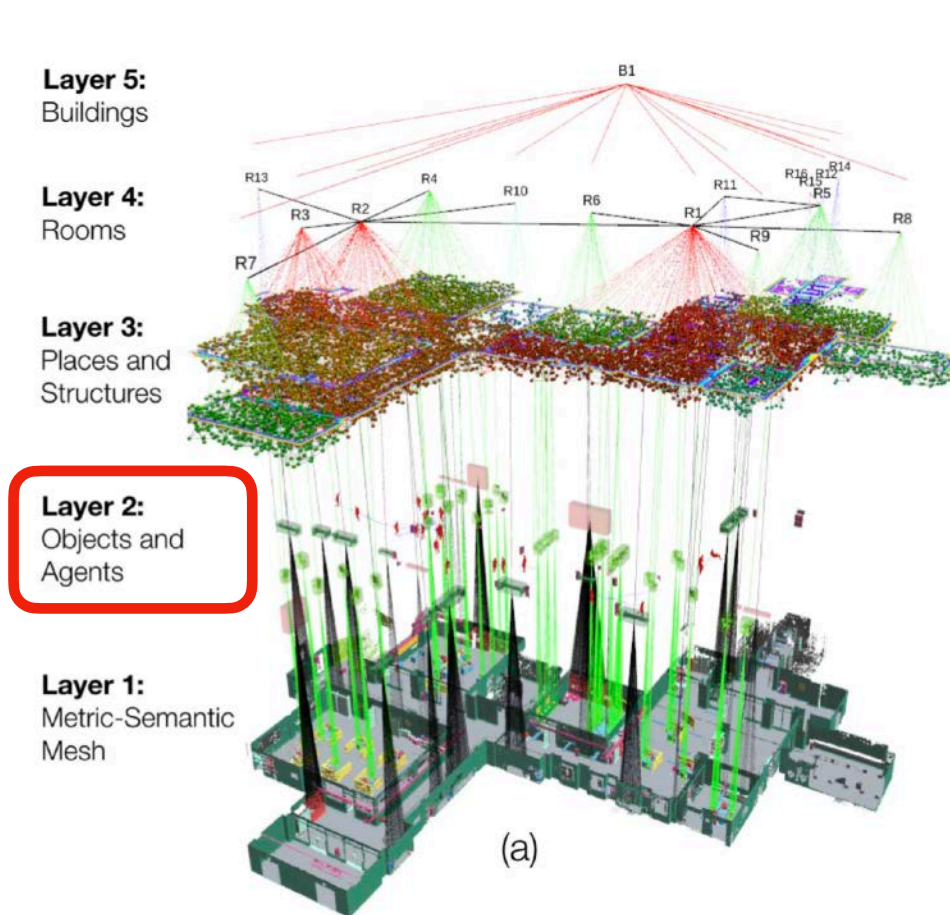


- From SLAM algorithms to a **Spatial Perception engine (SPIN)**, that infers geometry, semantics, a hierarchy of high-level spatial concepts and their relations

Figure 1 in Antoni Rosinol et al, "3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans." Robotics: Science and Systems 2020 Corvallis, Oregon, USA, July 12-16, 2020 © Antoni Rosinol et al. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

[Rosinol, Gupta, Abate, Shi, Carlone, 3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans, RSS'20]

# Layer 2: Objects and Agents



Our SPIN detects and tracks dense human models and builds a pose graph for further optimization and outlier rejection

## - Humans:

- 3D dense shape reconstruction from monocular images [2]
- Robust Pose Graph Optimization to track human poses over time

## - Objects:

- Euclidean clustering (when shape is unknown)
- TEASER++ (when shape is known)

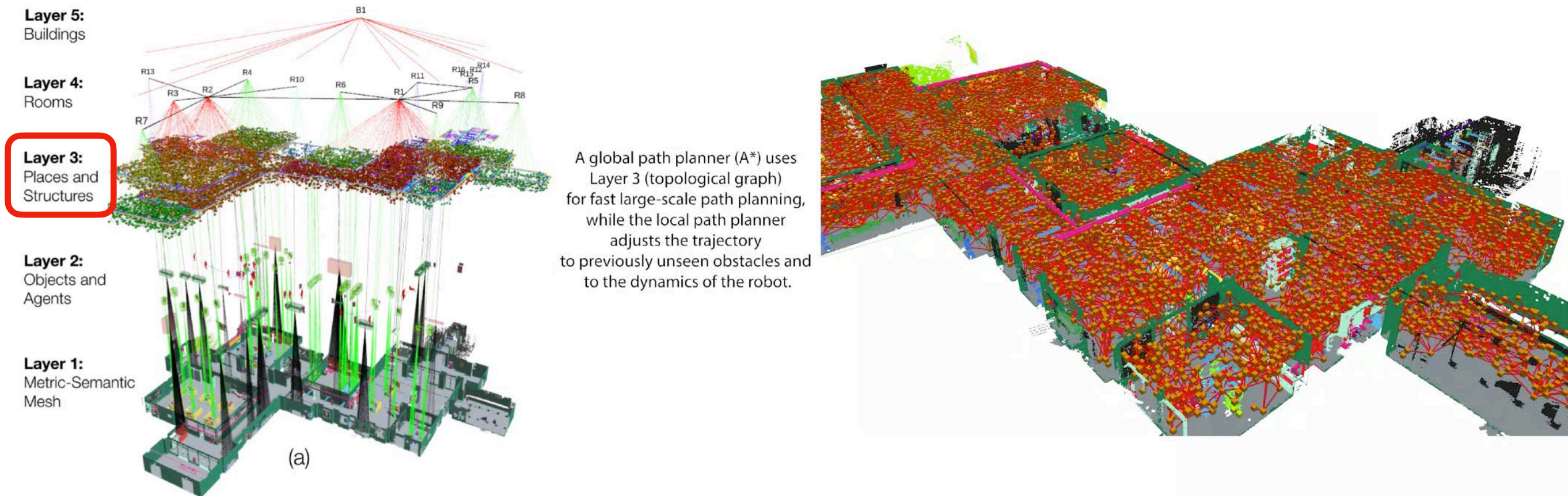
Figure 1 in Antoni Rosinol et al, "3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans." Robotics: Science and Systems 2020 Corvallis, Oregon, USA, July 12-16, 2020 © Antoni Rosinol et al. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

[1] Rosinol et al., 3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans, RSS'20.

[2] Kolotouros, Pavlakos, Daniilidis, Convolutional mesh regression for single-image human shape reconstruction, CVPR'19.



# Layer 3: Places and Structures



- **Places:** obstacle-free locations in the map, such that there is line-of-sight between pairs of nodes (suitable for fast path planning), using [2]
- **Structures:** separators between free space (walls, ground floor, ceiling)

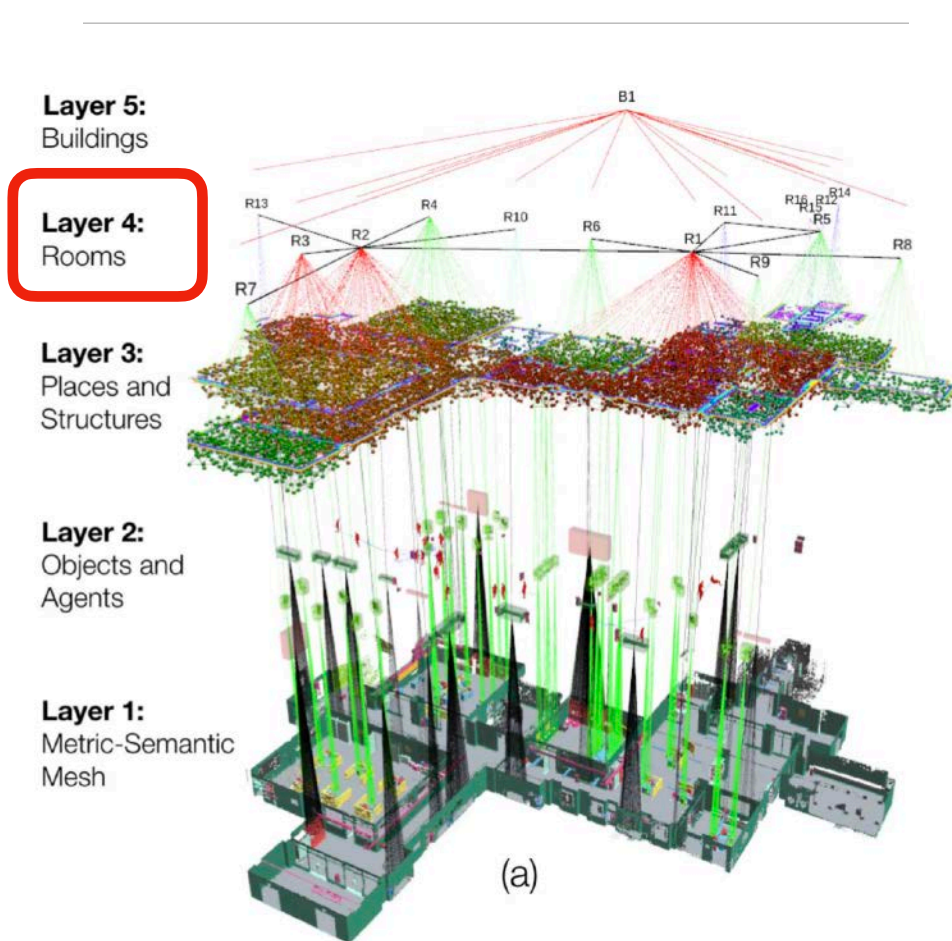
Figure 1 in Antoni Rosinol et al, "3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans." Robotics: Science and Systems 2020 Corvallis, Oregon, USA, July 12-16, 2020 © Antoni Rosinol et al. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <https://ocw.mit.edu/help/faq-fair-use/>

[1] Rosinol et al., 3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans, RSS'20. 38

[2] Oleynikova, Taylor, Siegwart, Nieto, Sparse 3D topological graphs for micro-aerial vehicle planning, IROS'18.



# Layer 4: Rooms



## Places and Room Clustering



We cluster the places in the environment into different rooms, obtaining an actionable representation for navigation and planning

## - Rooms:

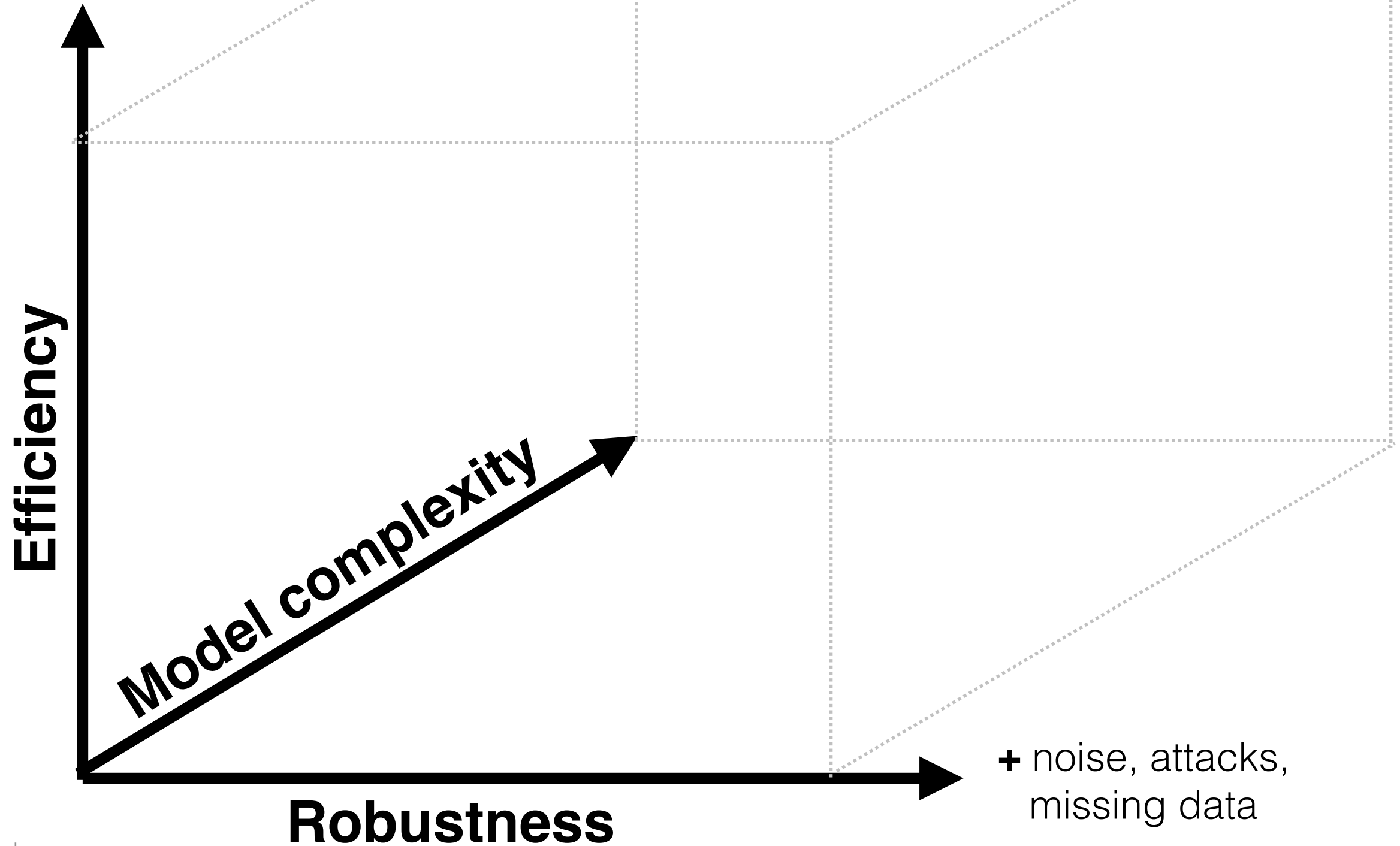
- extracted from graph of places using graph clustering

- **Remark:** traversability described at the level of rooms, places, and in the mesh: this is a “feature”, rather than a “bug” (-> hierarchical planning)



# Active Research Directions

- power, size,  
time constants



MIT OpenCourseWare  
<https://ocw.mit.edu/>

16.485 Visual Navigation for Autonomous Vehicles (VNAV)  
Fall 2020

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.