

Context for vehicles and sensing

- Environments:
 - Case 1:
 - work anywhere
 - likely: various gadgets improve safety and experience
 - implausible: full autonomy
 - Case 1.5:
 - work some places
 - some are specialized to freeways, etc.
 - Case 2:
 - work only in tightly controlled environment (eg smart city)
 - there are models of full autonomy (eg transporters at airports)

Case 1

- Various gadgets improve safety and experience
- case-by-case reasoning about representation and sensing
- Issues:
 - what's worth doing?
 - what can be done easily?
 - how much sensing?

Case 1 examples

- Reversing cameras
- Reversing sonar
- Forward sonar for collision avoidance
- Active collision management
- Pedestrian detection
- Various safety cameras
 - driver attention
 - record events for dispute resolution
 - driver sobriety
- Smarter links to maps

Case 1.5:

- Mostly, more specialized gadgets, mostly for highways
 - lane following for highways
 - predicting highway turnoffs
 - speed control that's aware of cars in front
 - neat tricks to reduce traffic jams
- Issues
 - what's worth doing?
 - what can be done easily?
 - how much sensing?

Case 2: Strongly controlled environments

- Full autonomy quite plausible
 - depending on regulatory and environmental control
 - there are models of full autonomy (eg transporters at airports)
 - This case is valuable, and may be important
 - public transport -> apartment in high density living areas
- Issues:
 - how much control do you need?
 - what density of traffic can be sustained?
 - how do you ensure safe behavior if weird stuff happens?

The questions that will plague us

- What representation do we need?
- How much data do we need to make it?
 - and where do we get it?
- How do we know if it works

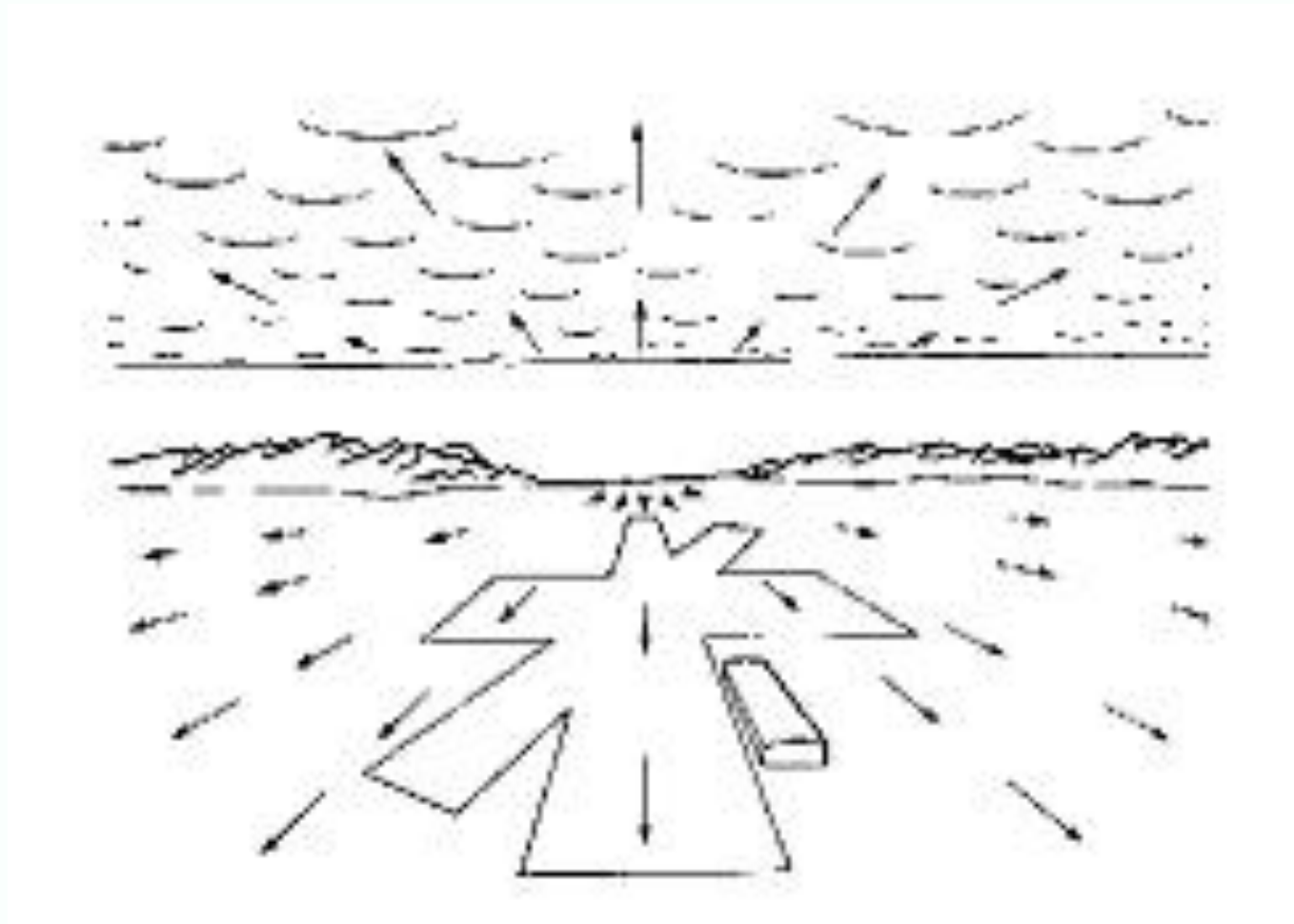
Representation

- A1: No representation required
 - link control inputs to sensing with multiple network layers
 - train on simulation with reinforcement learning
 - dubious position, but...
 - notice that, IN PRINCIPLE, this deals with full autonomy
 - Q:
 - how do you know it will do the right thing in a given situation?
 - A (dubious)
 - watch what it does on training data

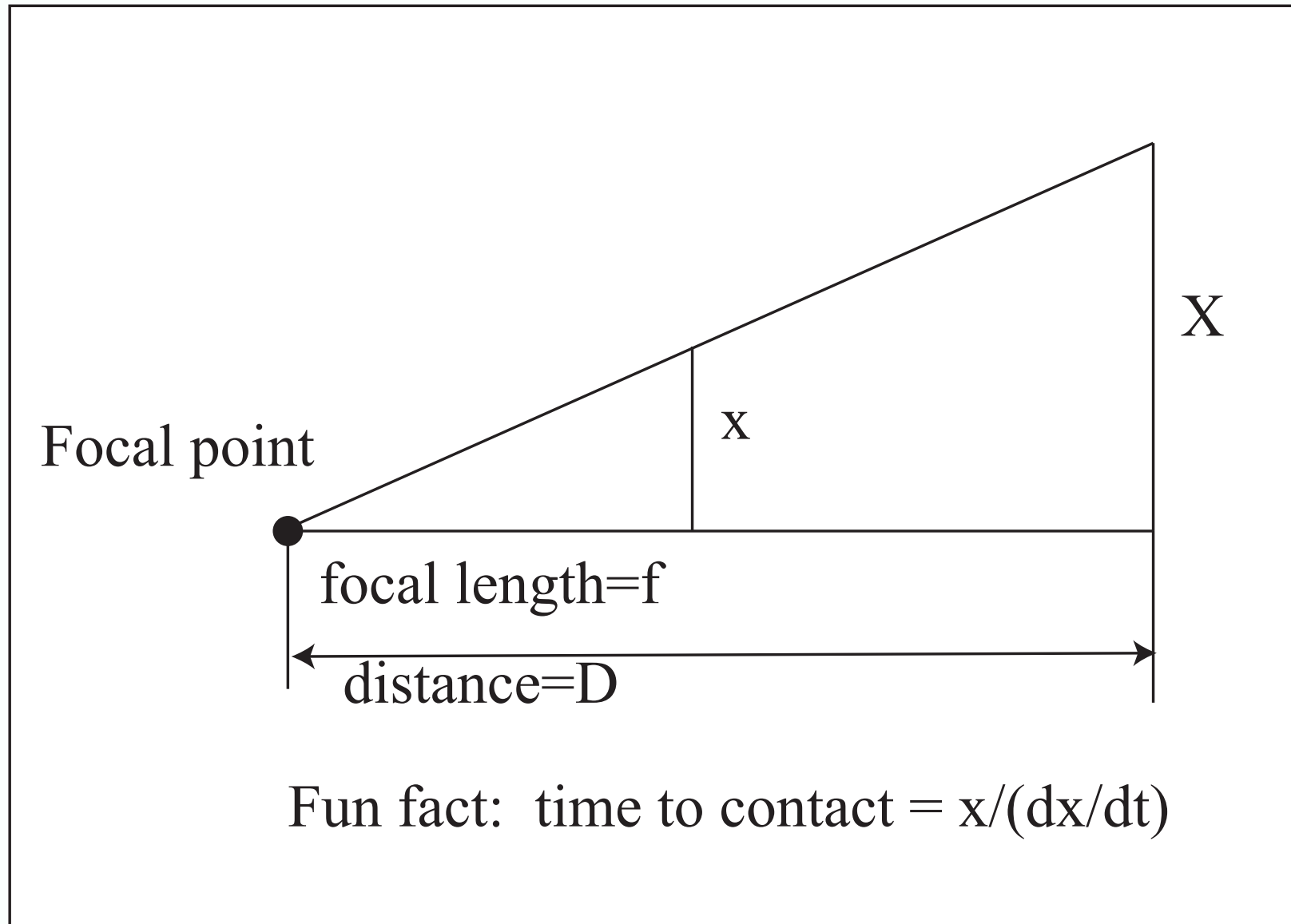
Representation

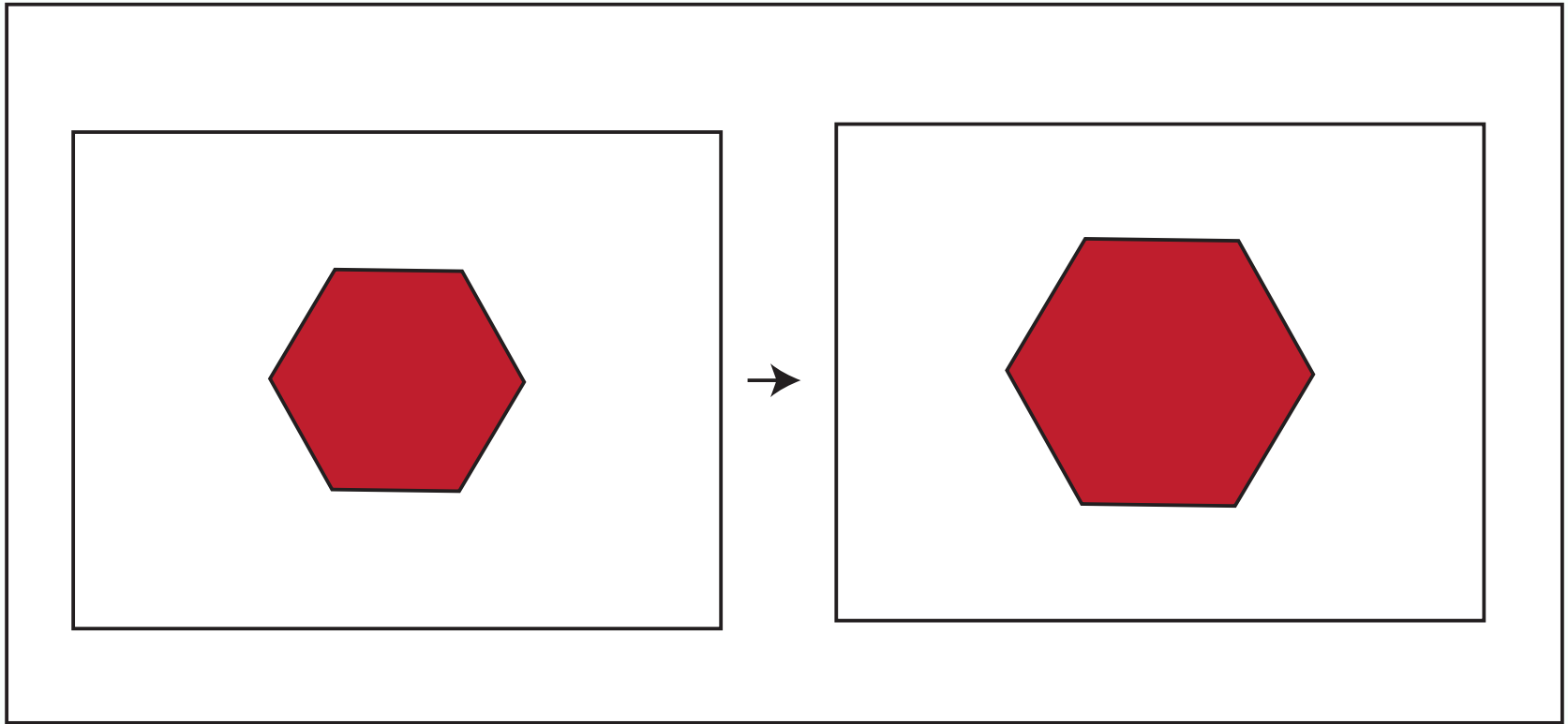
- A2: 3D reconstruction
 - build complete 3D model of world around you
 - LIDAR, SFM, etc.
 - label it with appropriate labels (next slides)
 - use a planner, etc to make paths in that environment
 - follow paths
 - Q:
 - how do you know it will do the right thing in a given situation?
 - A (dubious):
 - prove that environment is right and software is correct
 - Q:
 - do you really need a 3D representation?
 - A:
 - who knows?

Optic flow as a theory of perception

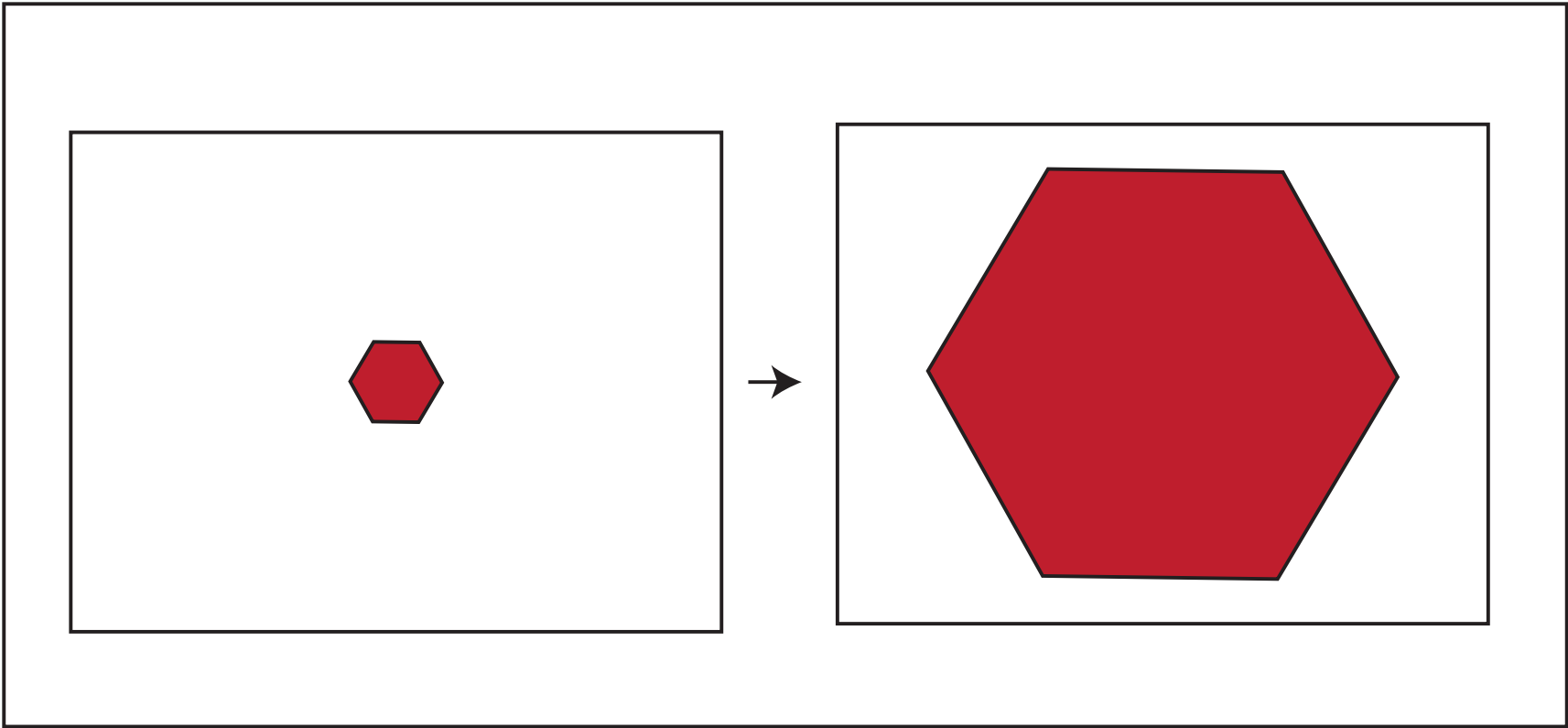


Fun fact about vision

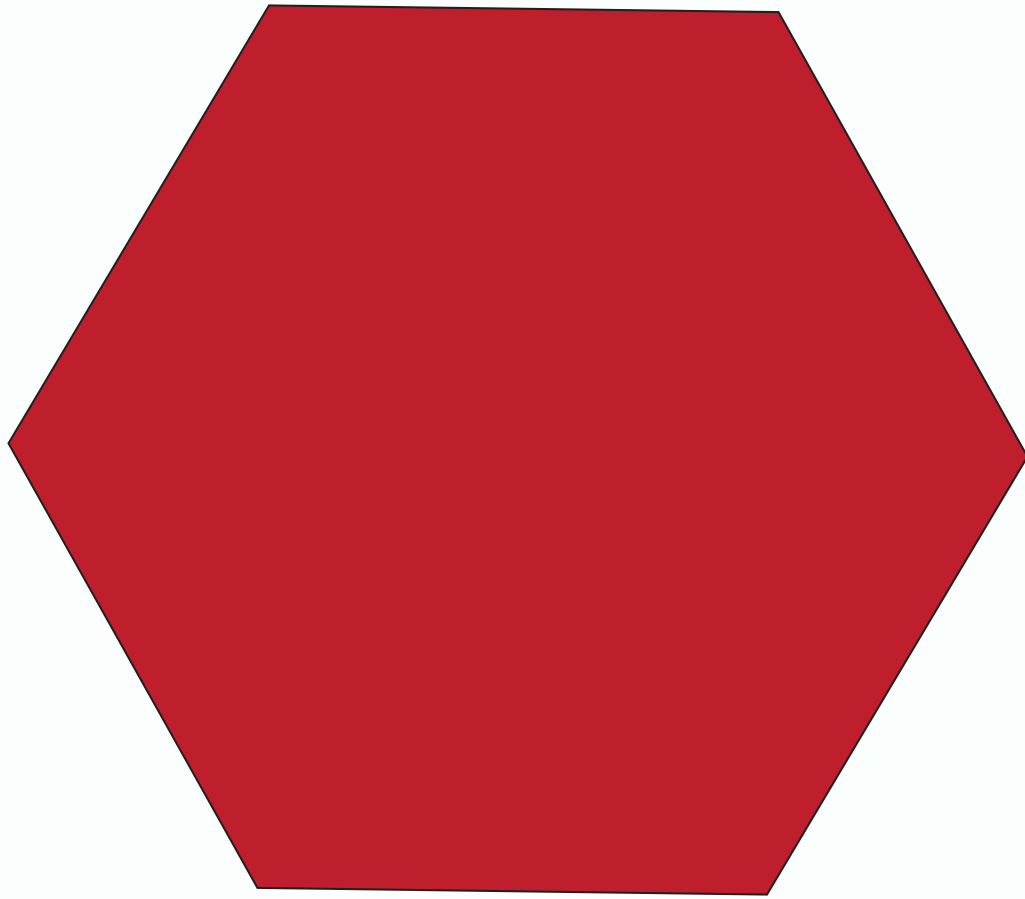


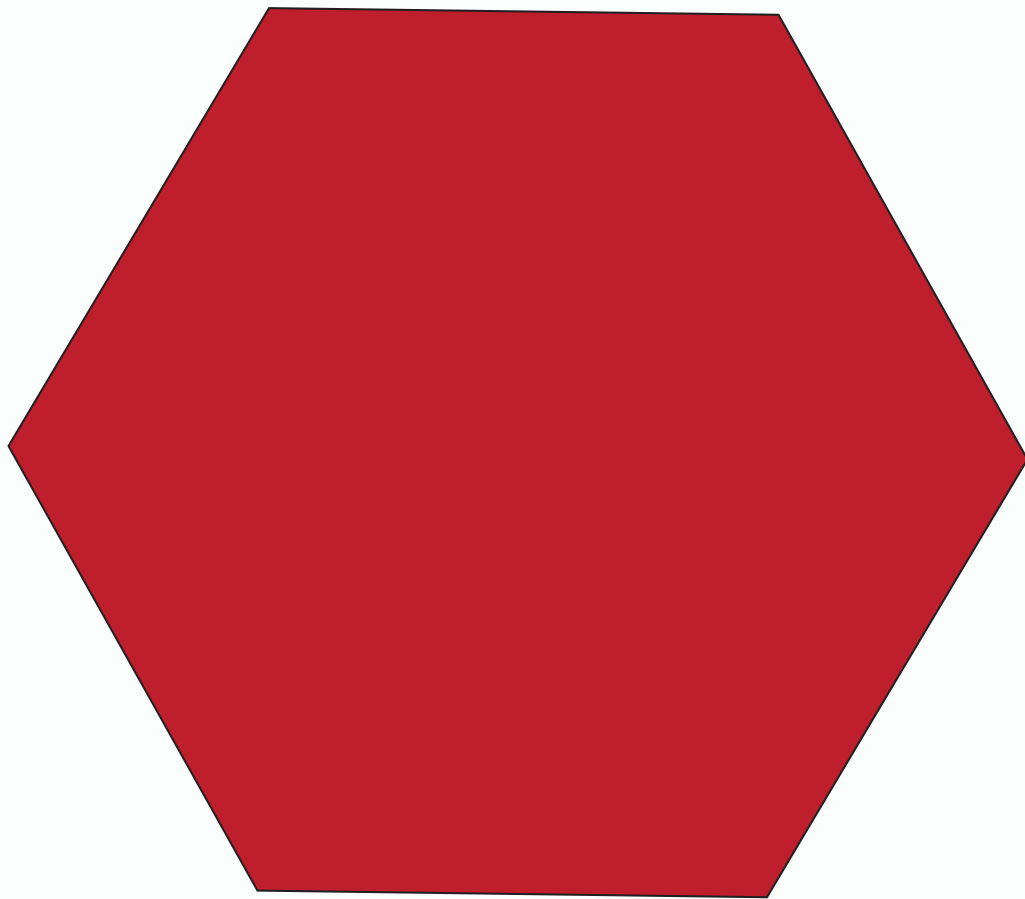


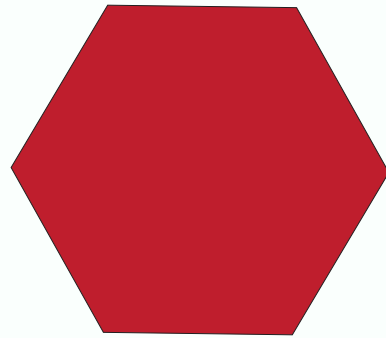
TTC - Long

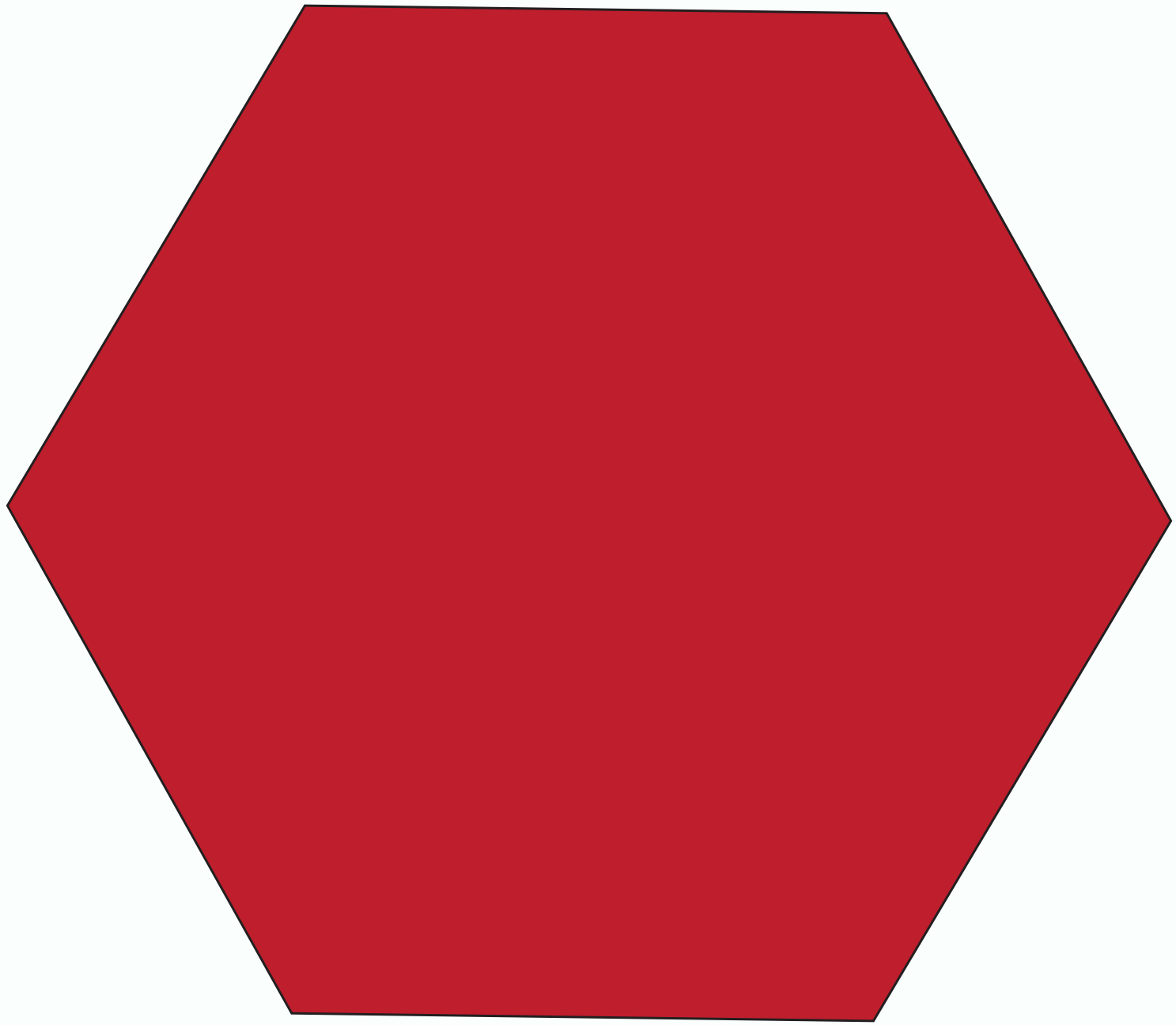


TTC - AAARGH!









Likely truth about 3D vs 2D

- Straightforward to convert from 2D to 3D repns and back
- This means anything you can do w/3D, you can do w/2D
- But: convenience is important
 - some planners want 3D
 - sensing 3D as 3D might be a good idea (LIDAR)
 - detection is generally faster in 2D, might be easier

Representation

- A3: Label images (or 3D reconstruction)
 - with what?
 - label all possible objects with all names
 - label some classes, ignore others
 - what taxonomy?
 - likely a derived taxonomy from actions
 - Q:
 - how do you know it will do the right thing in a given situation?
 - A (dubious):
 - prove that environment is right and software is correct
 - Q:
 - what should be labelled and what should be ignored?
 - A:
 - who knows? likely the things that most affect performance?

Labelling



Labelling



MS-CoCo

Labelling



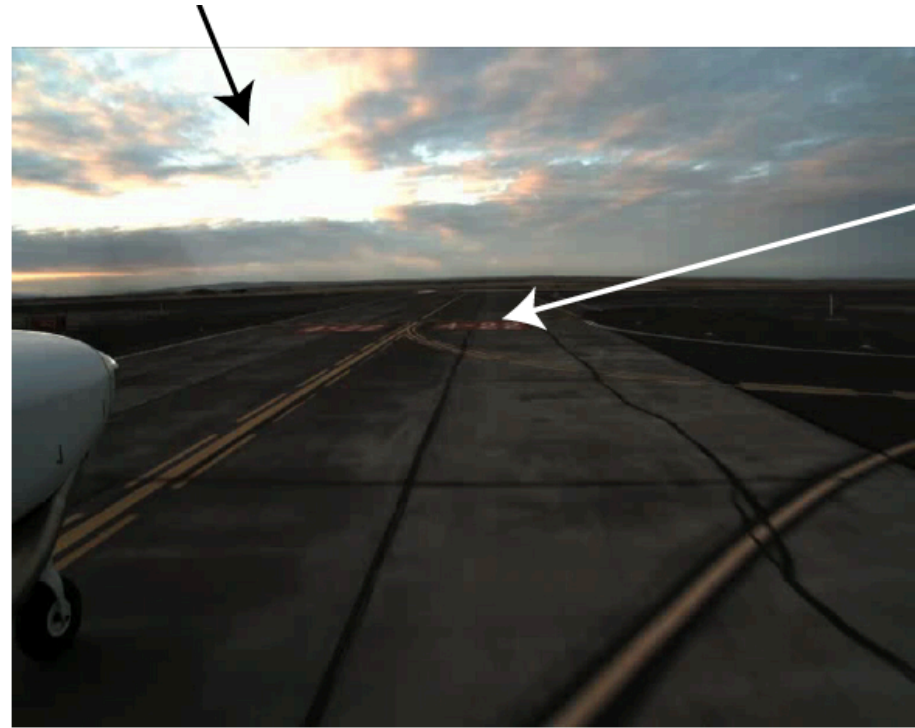
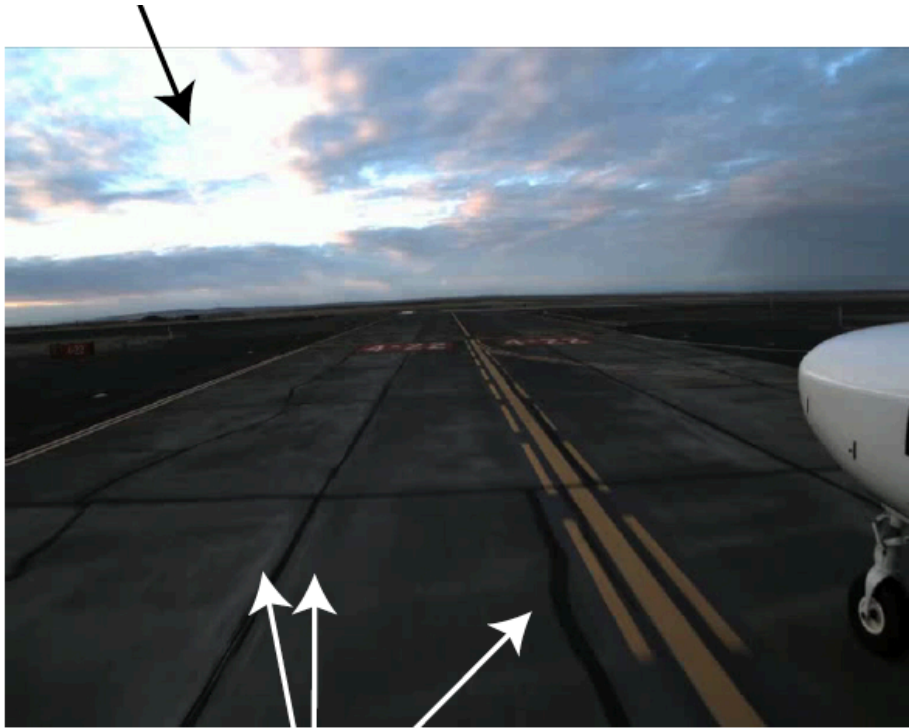
The questions that will plague us

- What representation do we need?
- How much data do we need to make it?
 - and where do we get it?
- How do we know if it works?

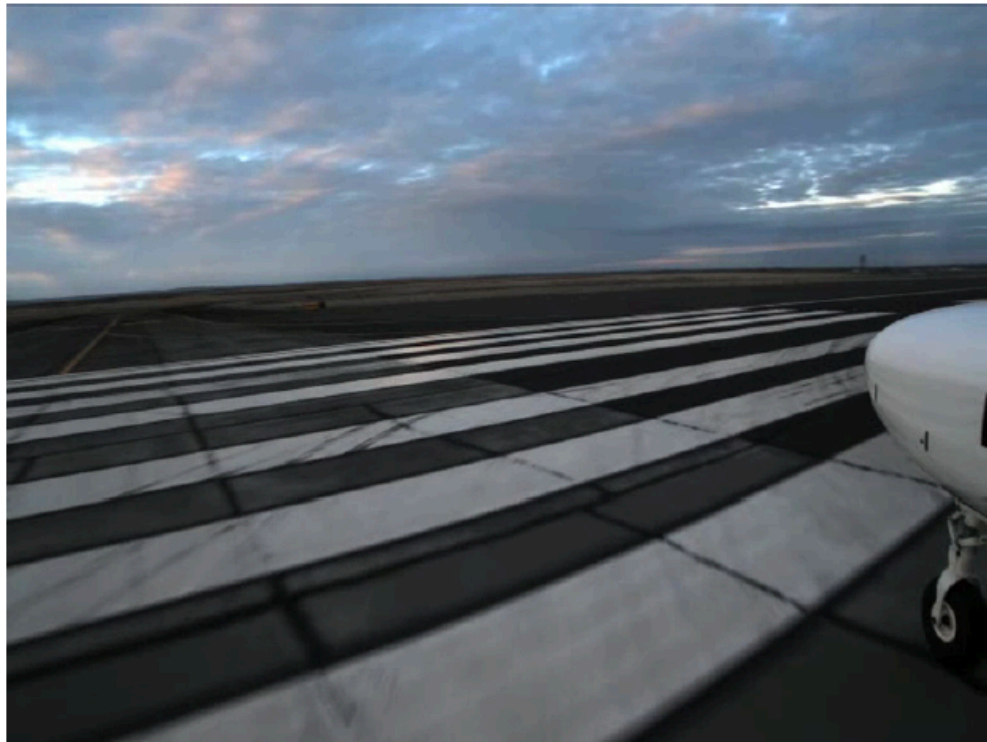
XXXX Autonomy data



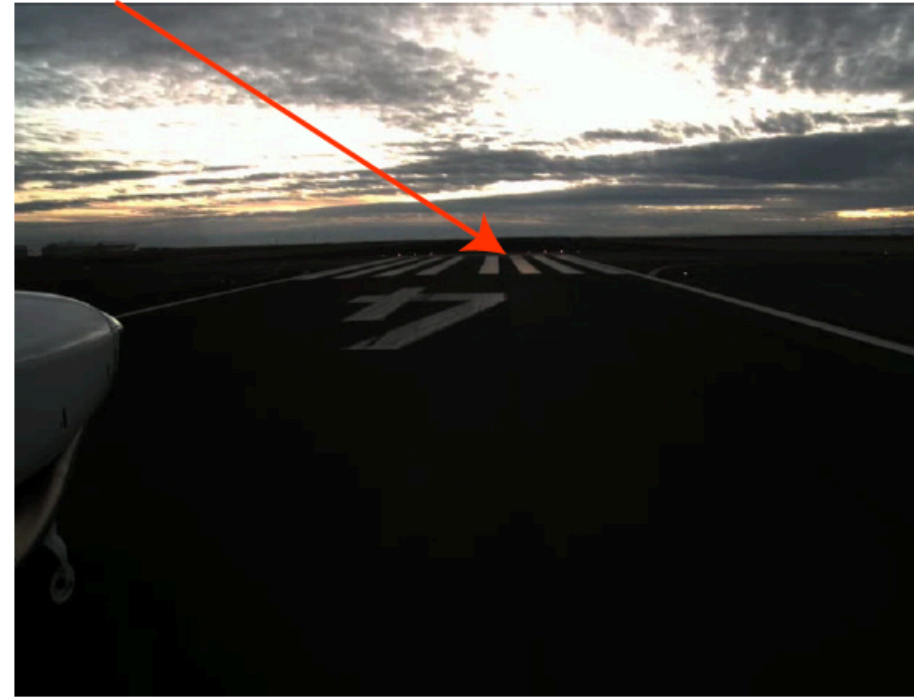
Special features: rich appearance variation



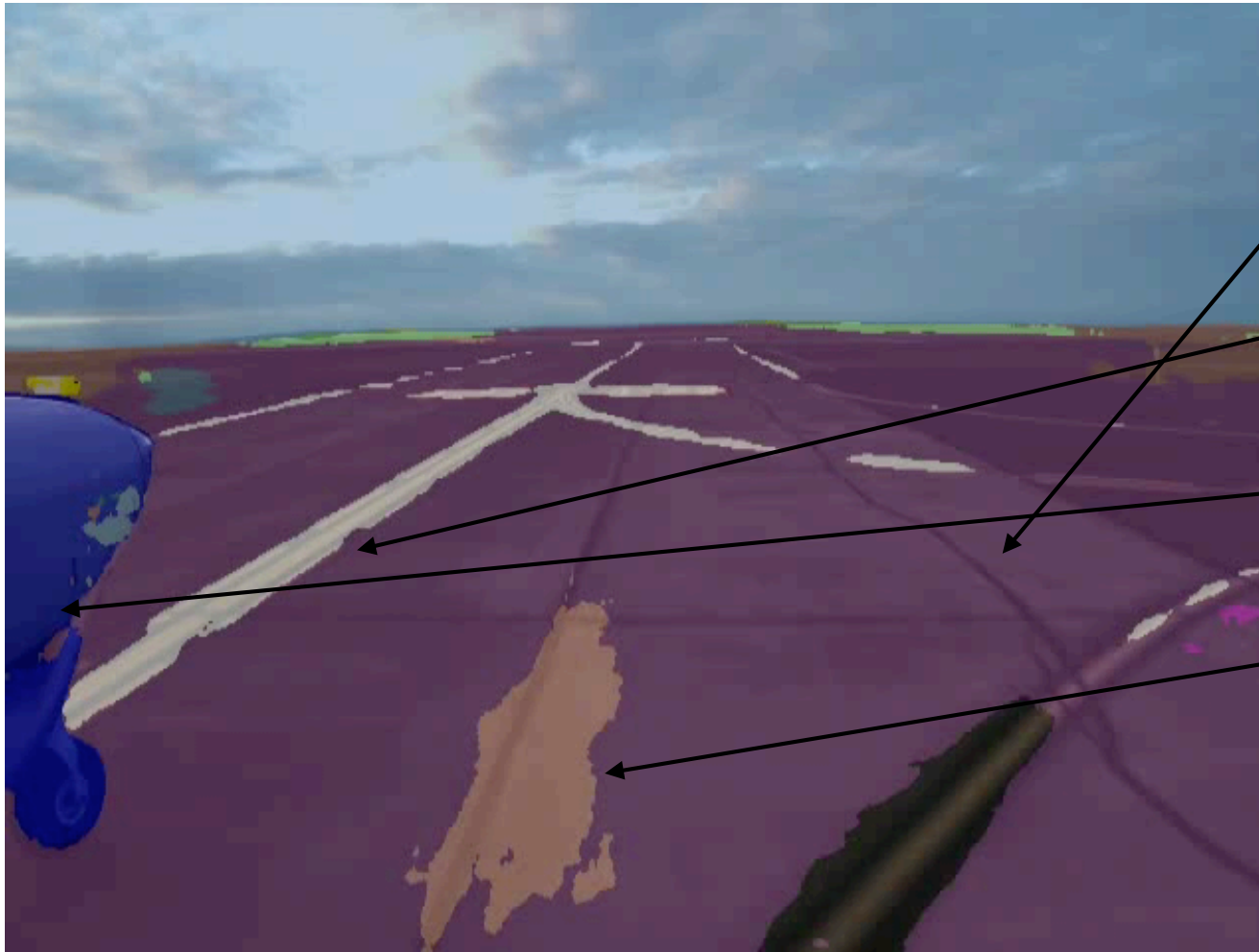
Special features: rich appearance variation



Special features: rich appearance variation



Standard semantic segmenter



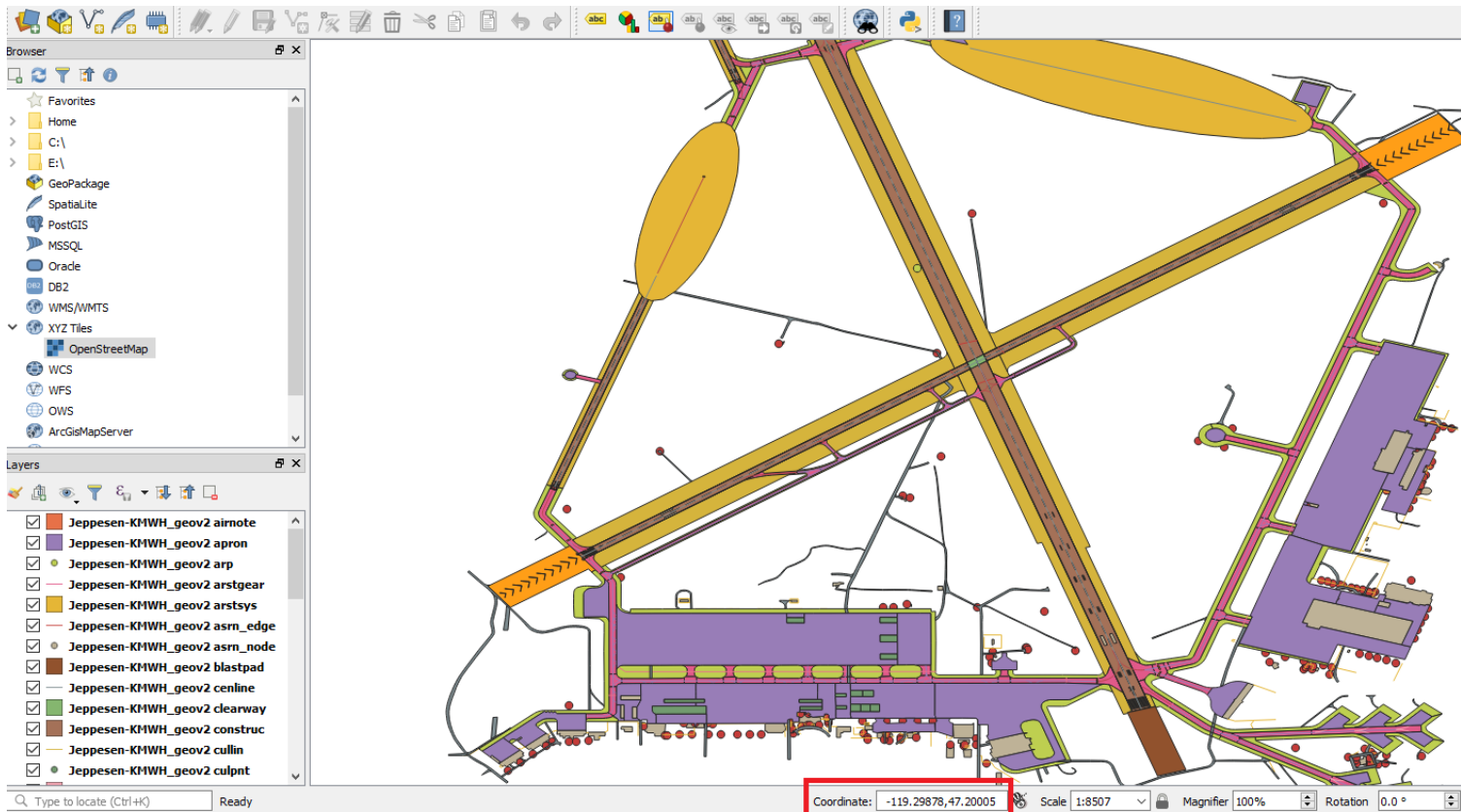
- Bird
- Ground Animal
- Curb
- Fence
- Guard Rail
- Barrier
- Wall
- Bike Lane
- Crosswalk - Plain
- Curb Cut
- Parking
- Pedestrian Area
- Rail Track
- Road
- Service Lane
- Sidewalk
- Bridge
- Building
- Tunnel
- Person
- Bicyclist
- Motorcyclist
- Other Rider
- Lane Marking - Crosswalk
- Lane Marking - General
- Mountain
- Sand
- Sky
- Snow
- Terrain
- Vegetation
- Water
- Banner
- Bench
- Bike Rack
- Billboard
- Catch Basin
- CTV Camera
- Fire Hydrant
- Junction Box
- Mailbox
- Manhole
- Phone Booth
- Pothole
- Street Light
- Pole
- Traffic Sign Frame
- Utility Pole
- Traffic Light
- Traffic Sign (Back)
- Traffic Sign (Front)
- Trash Can
- Bicycle
- Boat
- Bus
- Car
- Caravan
- Motorcycle
- On Rails
- Other Vehicle
- Trailer
- Truck
- Wheeled Slow
- Car Mount
- Ego Vehicle

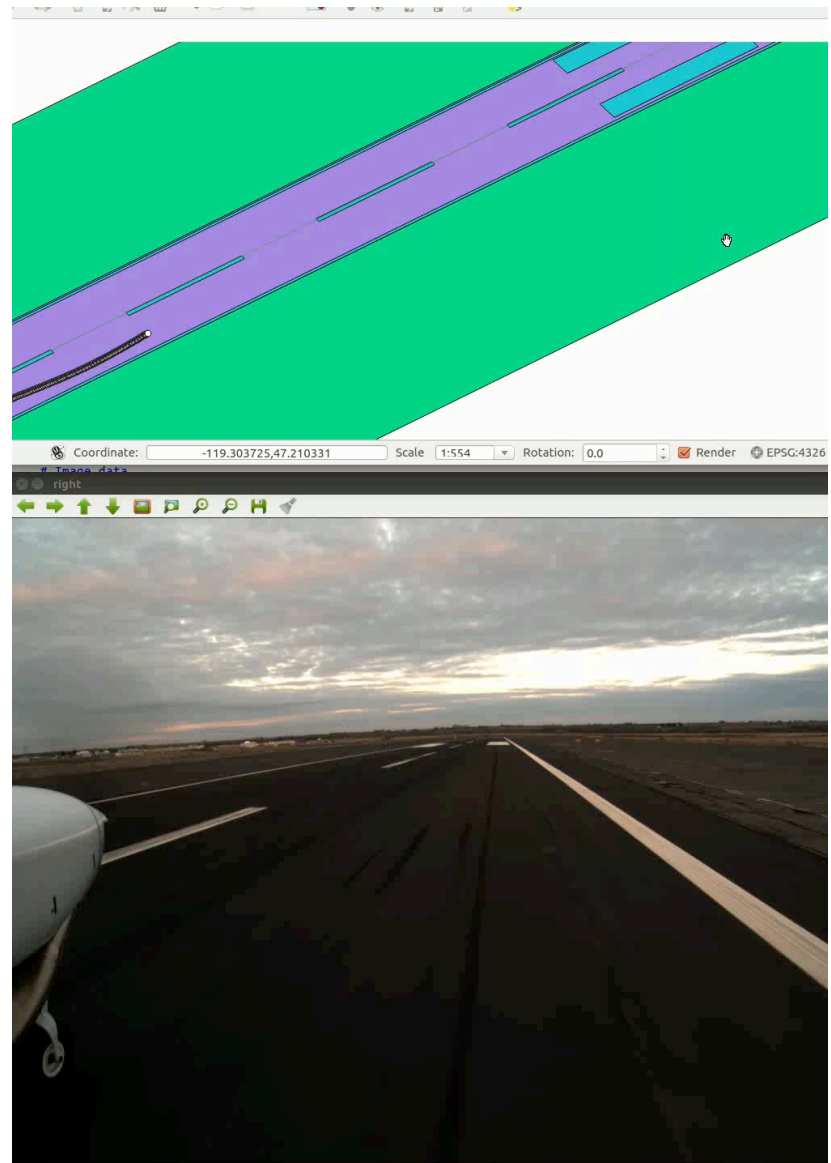
XXX data consequences



XXXX data consequences







GPS is quite good, but not perfect

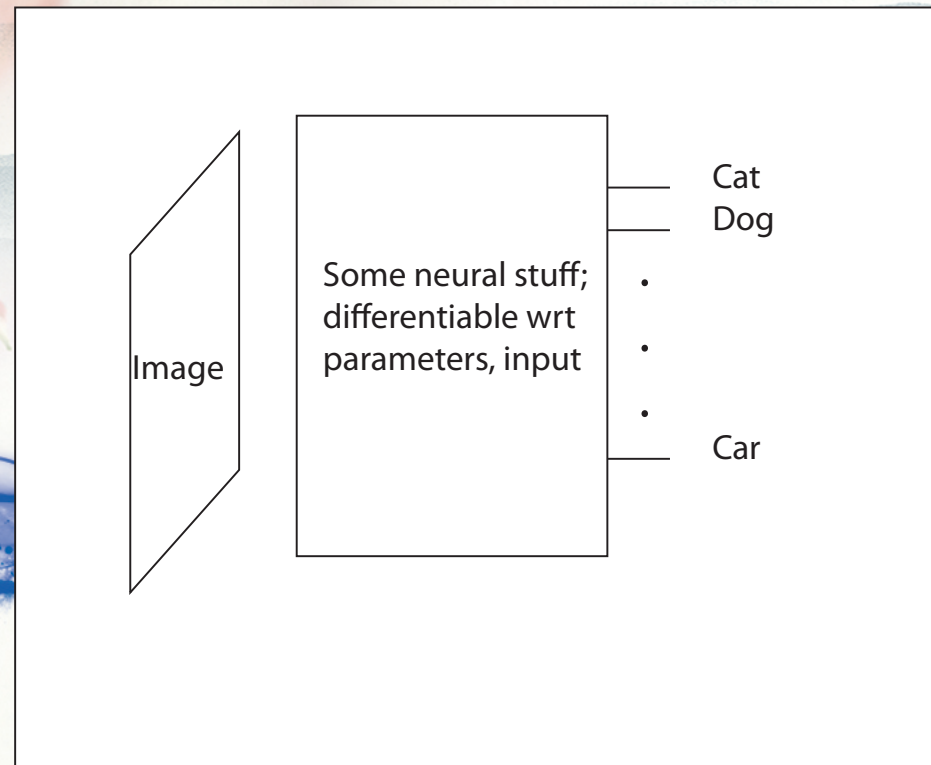
Another curious data problem



The questions that will plague us

- What representation do we need?
- How much data do we need to make it?
 - and where do we get it?
- How do we know if it works?

Image classification



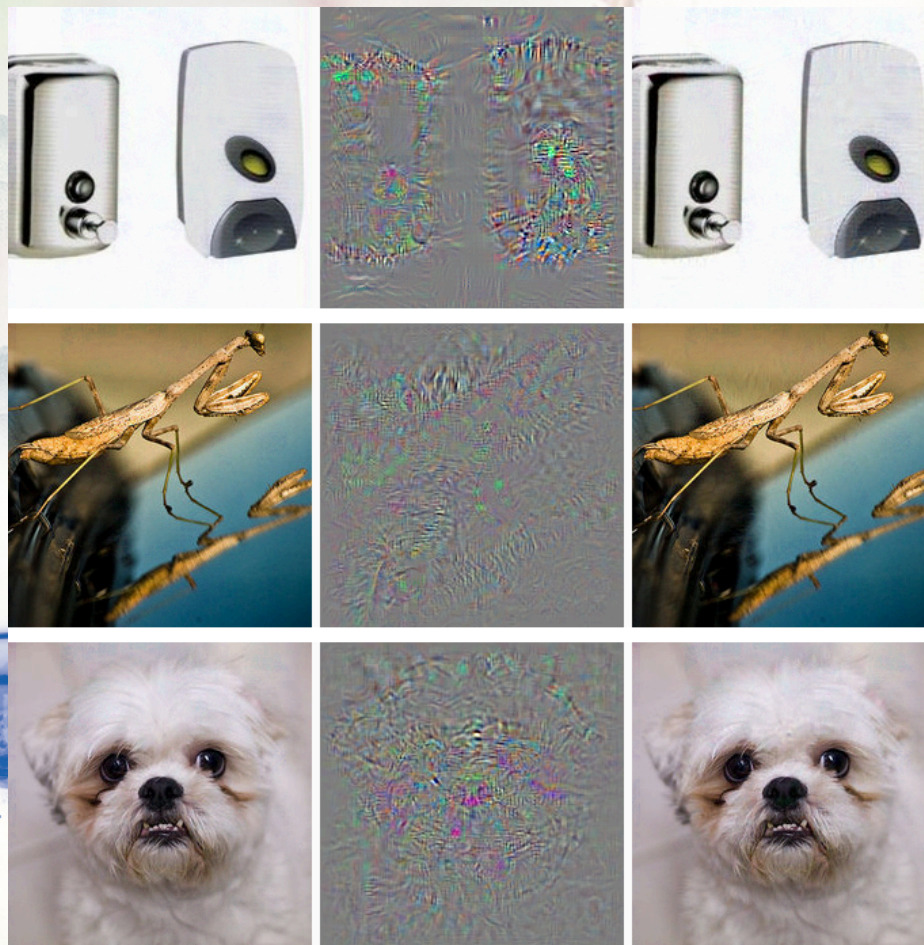
Real-Time
Perception/ Prediction
Behavior Forecast
computer vision technology
Real-Time Perception
Zahar
Object Detection
Behavior Forecasting
Next Gen Simulation
oard
ration

Adversarial example

- Search for
 - small update to image
 - such that
 - output for true class is low
 - output for some other class is high
- Surprising fact:
 - such updates can be **VERY** small



Correctly
classified

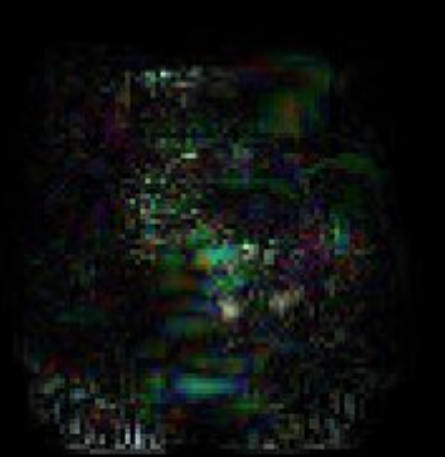
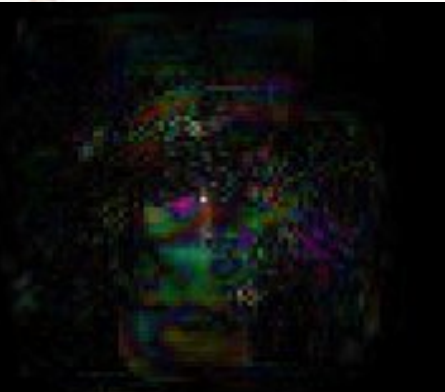


APOLLOSCAPE

“Ostrich”

Szegedy et al, 13

Real-Time Perception/Prediction
Behavior Forecasting
Computer vision technology
Heat-Time Perception
Localization
Object Detection
Behavior Forecasting
Next-Gen Simulation
On-board
Camera calibration



APOLLOSCAPE

Szegedy et al, 13

Real-Time Perception/Prediction Behavior Forecasting
Computer vision technology
Real-Time Perception/Prediction Behavior Forecasting
Object Detection Behavior Forecasting
Next Gen Simulation
board

camera calibration

Fast gradient sign search

- Search sign(gradient)

- x is image

- J is some cost

- eg J=(true class)-(best false class)

- Iterate

$$\eta = \epsilon \text{sign}(\nabla_x J(\theta, x, y)).$$

$$\mathbf{X}_{adv} = \mathbf{X} - \epsilon \cdot \text{sign}(\nabla_{\mathbf{X}} J(\mathbf{X}, y_{fool})).$$

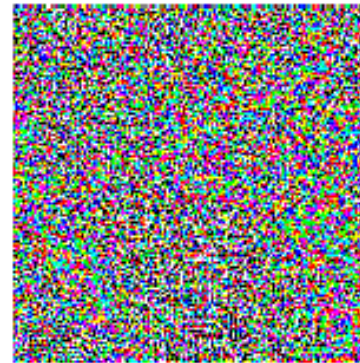


Fast gradient sign



x
 “panda”
 57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$
 “nematode”
 8.2% confidence

=



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
 “gibbon”
 99.3 % confidence



Deepfool

- Find r by
 - linearize k
 - update
 - (possibly) repeat

$$\Delta(x; \hat{k}) := \min_r \|r\|_2 \text{ subject to } \hat{k}(x+r) \neq \hat{k}(x)$$

label

image

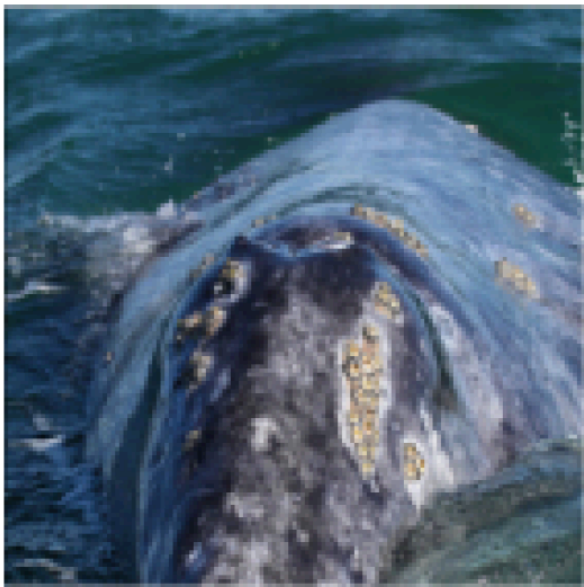
update



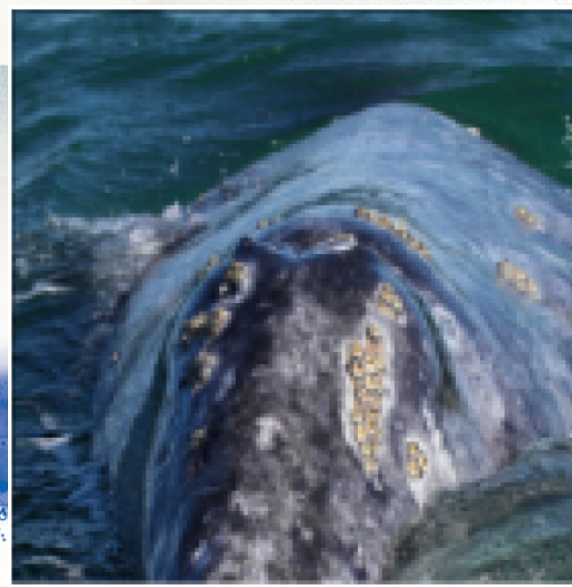
Moosavi-Desfooli et al 16

Deepfool

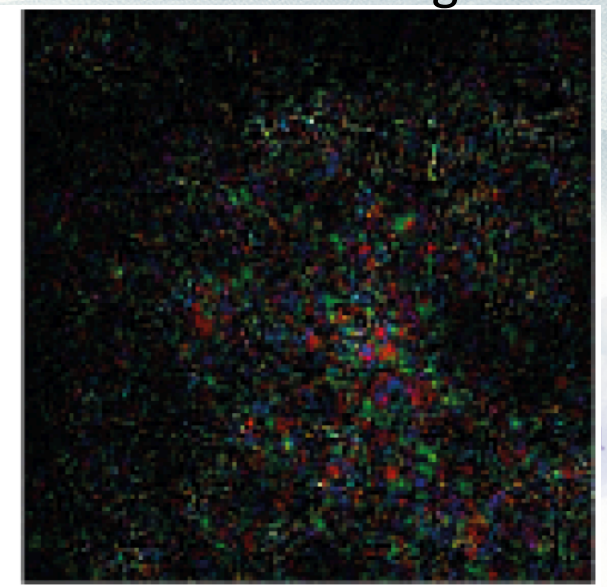
Whale



Turtle



Difference image

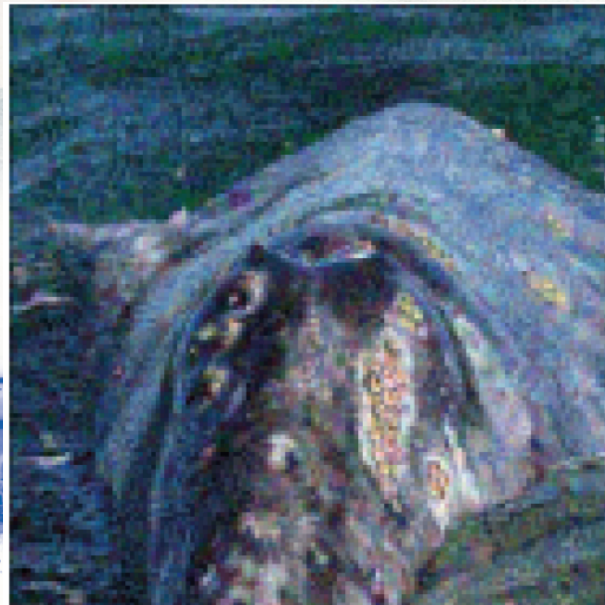


Compare fast gradient sign

APOLLOSCAPE

Turtle

Difference image



Moosavi-Desfooli et al 16

Flow based methods

- New image obtained by

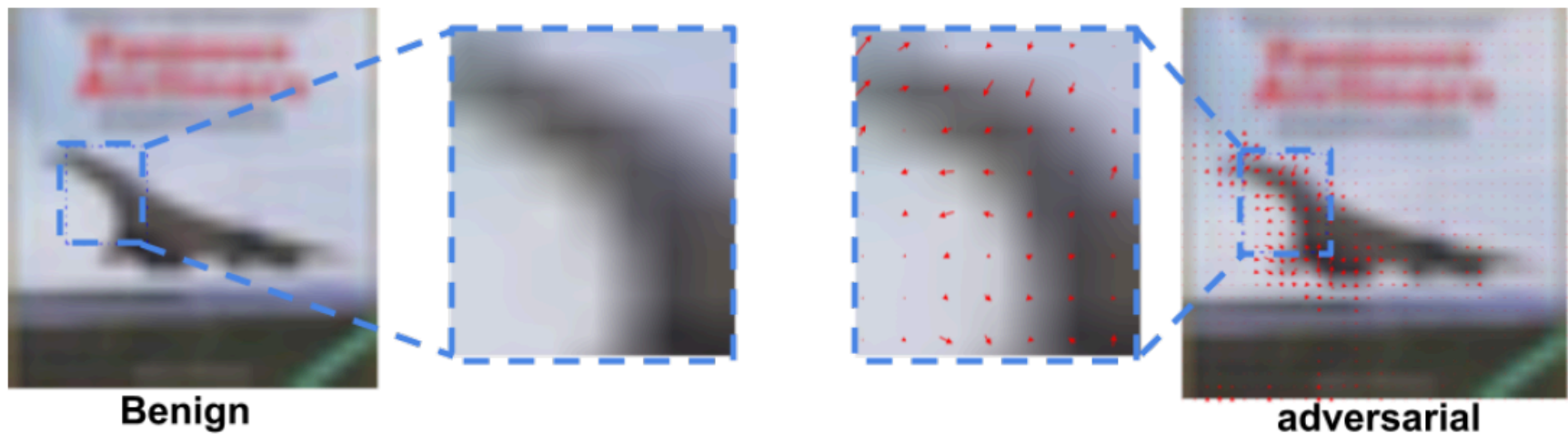


Figure 6: Flow visualization on CIFAR-10. The example is misclassified as bird.

Xiao et al
2018

Flow based methods

- New image obtained by moving pixels

$$\mathcal{X}_{adv}(u, v) = \mathcal{X}(u + f_u, v + f_v)$$

$$f^* = \operatorname{argmin}_f \mathcal{L}_{adv}(x, f) + \tau \mathcal{L}_{flow}(f)$$

Size of flow

Adversarial loss



Flow methods

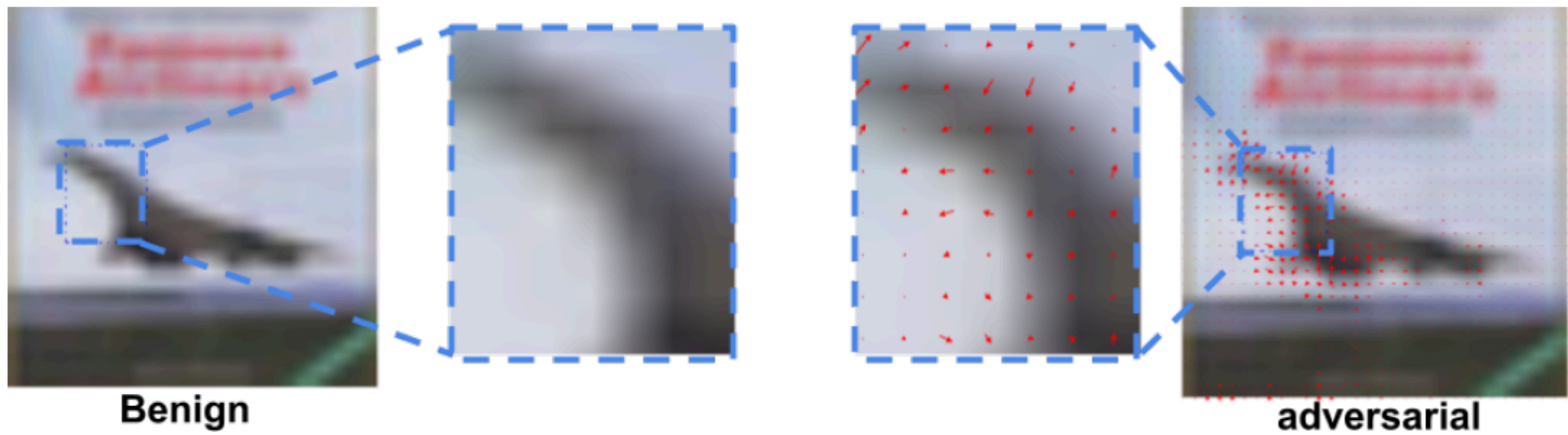


Figure 6: Flow visualization on CIFAR-10. The example is misclassified as bird.

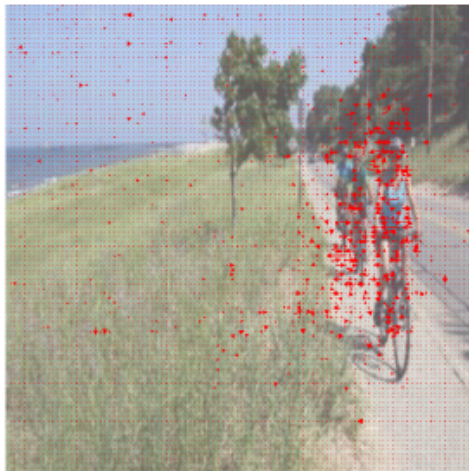
Xiao et al 2018

Flow methods

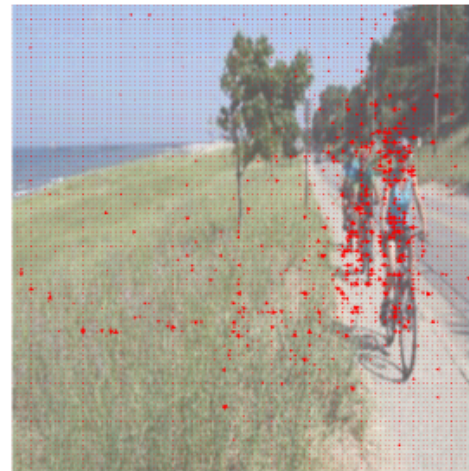
APOLLOSCAPE



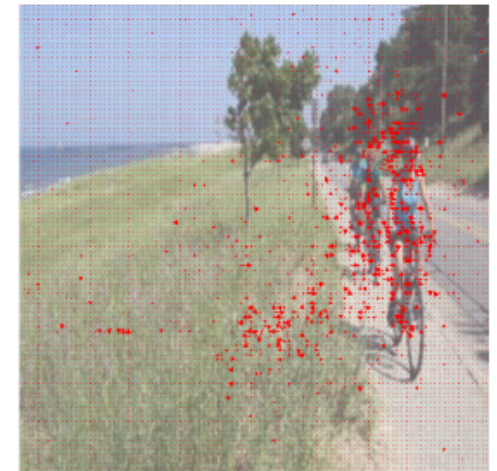
(a) mountain bike



(b) goldfish



(c) Maltese dog



(d) tabby cat

interfaces On-board Behavior Forecasting
camera calibration Next Gen Simulation

Xiao et al 2018



APOLLOSCAPE

Real-Time
Perception/ Prediction
Behavior Forecasting
Computer vision technology
Object Detection
Behavior Forecasting
Next Gen Simulation

Lu et al 18

Yolo attack

- Yolo uses a large image area to
 - predict boxes
 - predict classes
- This means that a detection is
 - affected by pixels OUTSIDE box





