# EPFL VITA

**Visual Intelligence for Transportation**

"Humans subconsciously **forecast the future**…
Autonomous Vehicles must have the same **forecasting** capability to **harmlessly** and **effectively co-exist**",
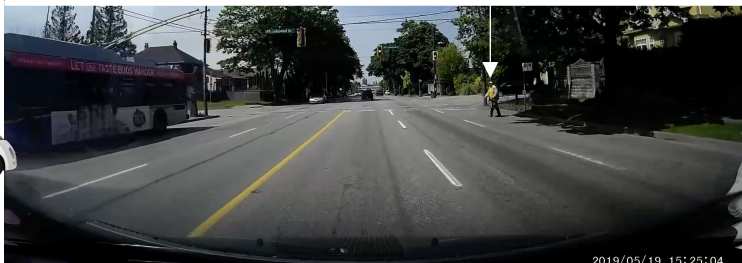Our lab goal.

**Forecasting** is essential…
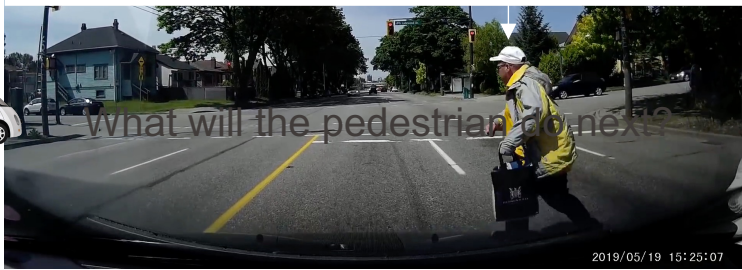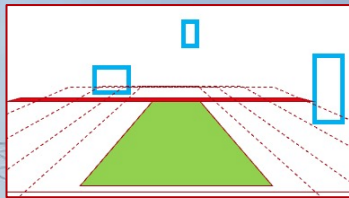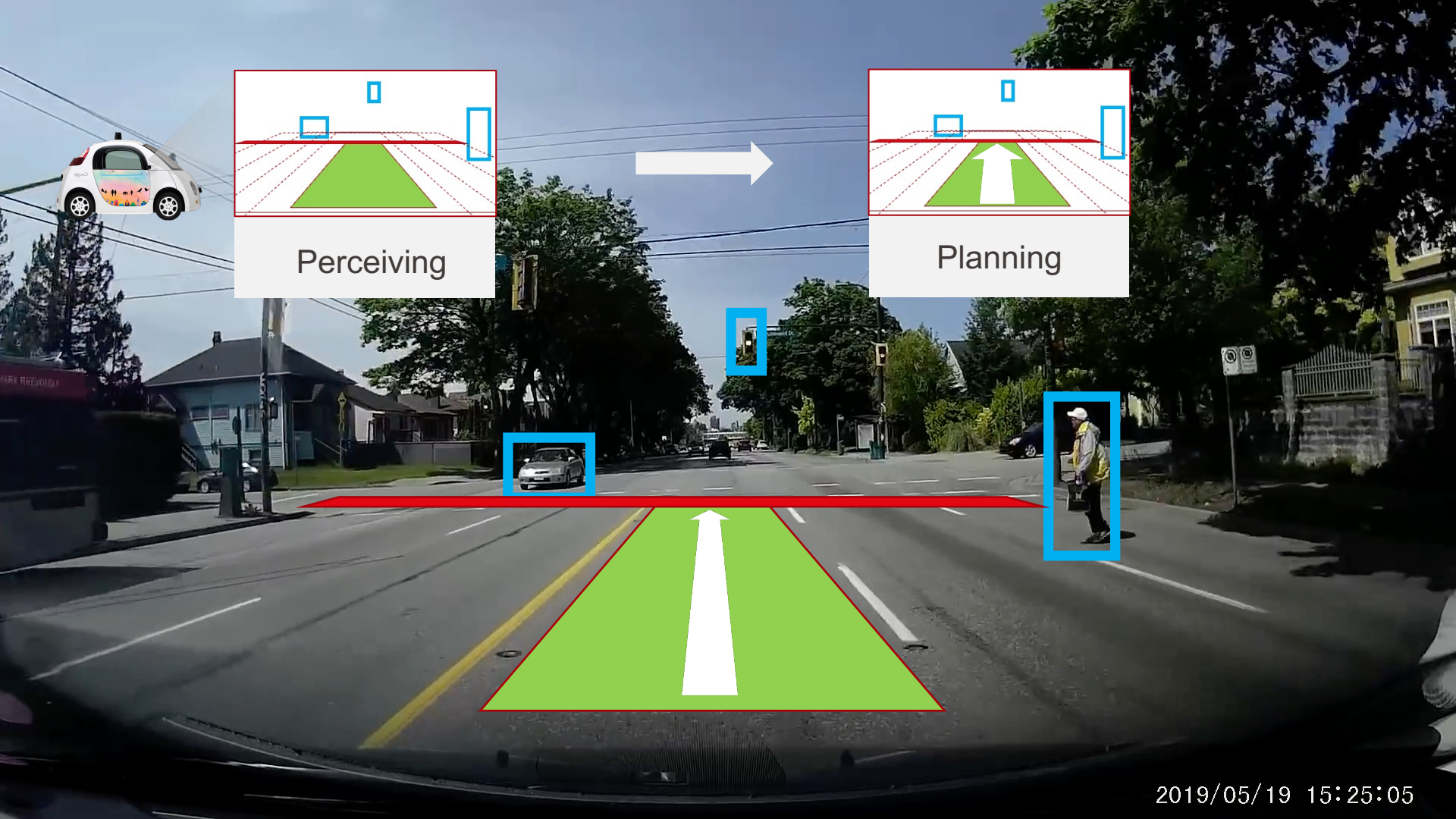
time   t

Autonomous ✓

t+1

Autonomous ✓

t+2

t+3

What will the pedestrian do next?

Perceiving

Planning

2019/05/19 15:25:05

Perceiving

Forecasting

Planning

Forecast to stop

2019/05/19 15:25:05

AI must forecast agent-agent* **interactions** =   **Social Forecasting**

*agent = any moving entity in the world (driver, pedestrian, cyclist…)

Forecast **not** to stop

**98% of AV accidents are due to an unexpected stop***

* Favarò,F., et al., "Examining accident reports involving autonomous vehicles." PLoS one , '17

Autonomous ✓

Autonomous ✓

Autonomous ✗
**=> Our robot freezes in close human proximity**

Autonomous ✗
**=> Our robot does not comply with social norms**

time t

t+1

t+2

t+3

VITA

Social Forecasting



Perceiving

**Socially-aware AI**

Planning

**Object detection**

# Body poses + Activities + Relationships

Standing   Standing

Walking

Walking

Talking

2.3m        2.1m    2.2m    4.5m

[1] PifPaf: Composite Fields for Human Pose Estimation, **CVPR'19; on-line demo:** vitademo.epfl.ch/movements

[2] Keypoints communities, **ICCV'21;** [3] OpenPifPaf**, open-source library, IEEE ITS'21**

[4] Convolutional Relational Machine, **CVPR'19;** [5] Detecting 32 Pedestrian Attributes, **IEEE ITS'21**

[6] Monoloco: Monocular 3D pedestrian localization and uncertainty estimation, **ICCV'19;** [7] MonStereo, **ICRA'21**

[1] Monoloco: Monocular 3D pedestrian localization and uncertainty estimation, **ICCV'19;** [2] MonStereo, **ICRA'21**

[1] Monoloco: Monocular 3D pedestrian localization and uncertainty estimation, **ICCV'19;** [2] MonStereo, **ICRA'21**

Object detection then Object tracking

# Pose detection with tracking

[1] OpenPifPaf, **IEEE transactions on ITS'21**          (online demo: https://vitademo.epfl.ch/movements)

[1] OpenPifPaf, **IEEE transactions on ITS'21**        (online demo: https://vitademo.epfl.ch/movements)

Body text:

[1] **on-line demo:** vitademo.epfl.ch/movements
[2] Keypoints communities, **ICCV'21;** [3] OpenPifPaf, open-source library, **IEEE transactions on ITS'21**

Orange box:
adult
female
walking laterally
alone
looking
intention to cross
bag on the right side
clothes lower dark
clothes upper light
pose right
has a stroller cart

Yellow box:
adult
female
walking laterally
alone
not looking
no intention to cross
no bag
clothes lower dark
clothes upper dark
pose front

Blue box (top):
adult
female
walking along the road
alone
not looking
no intention to cross
no bag
clothes lower dark
clothes upper dark
pose back

Green box:
adult
male
walking laterally
alone
looking
intention to cross
no bag
clothes lower dark
clothes upper dark
pose front

Blue box (bottom):
adult
female
walking laterally
alone
not looking
no intention to cross
no bag
clothes lower dark
clothes upper dark
pose back

[1] Detecting 32 Pedestrian Attributes for Autonomous Vehicles, **IEEE transactions on ITS'21**

[1] Looking dataset: https://looking-vita-epfl.github.io/

Intention to cross (in blue), or not (in green)

[1] Pedestrian Intention Prediction: A Convolutional Bottom-Up Multi-task Approach, **TRC'21**

Alexandre Alahi, EPFL

VITA

Social Forecasting

Perceiving

**Socially-aware AI**

Planning

# Social Forecasting
## (w/ pedestrians)

- **Input**: several sequences of states

- **Output**: forecast the future states, *e.g.*, next 5 seconds



Alexandre Alahi, EPFL

Observed sequence
Forecasted sequence

$(x_1^1, y_1^1)$  $(x_1^2, y_1^2)$  $(x_1^3, y_1^3)$

VITA

# Social Forecasting
## (w/ pedestrians)

- **Input**: several sequences of states

- **Output**: forecast the future states,
  *e.g.*, next 5 seconds

- **State**:
  - $(x^t, y^t)$ coordinates in time
  - Body pose [1]
  - Attributes (*e.g.*, on the phone, eye contact) [2]

- Challenge 1: agent-agent interactions
- Challenge 2: disentangle physics from social

[1] PifPaf, CVPR'19
[2] 32 attributes detector, ITS transactions'21



On the phone

Eye contact

Observed sequence
Forecasted sequence

# Social Forecasting
## (w/ vehicles)

Alexandre Alahi, EPFL

- **Input**: several sequences of states + scene infrastructure

- **Output**: forecast the future states, *e.g.*, next 5 seconds

- Challenge 1: agent-agent interactions

- Challenge 2: agent-scene interactions

- Challenge 3: additional external constraints



$(x_1^1, y_1^1)$  $(x_1^2, y_1^2)$

→ Observed sequence
⇢ Forecasted sequence

# Social Forecasting
## (w/ pedestrians)

# Social Forecasting
## (w/ vehicles)

→ Observed sequence
⇢ Forecasted sequence

# Robustness

# Learning paradigms

We collected the **largest** crowd dataset
(42 millions of examples)

✗ **Accuracy**
✓ **Interpretability**
✓ **Robustness**

2016
multi-class SF
(ECCV)

**Our Work**

Knowledge

Other Work

1995
Social Force
(Physical review)

1995-2016
Discrete Choice Models
RVO, ORCA, IGP

VITA

SF = Social Force
IGP = Interacting Gaussian Processes
RVO = Reciprocal Velocity Obstacle
ORCA = Optimal reciprocal collision-avoidance
LSTM = Long Short-Term Memory
GAN = Generative Adversarial Network
TTT = Test-Time Training

# Trajnet++

- Open-source library (> 15 models)
  - https://github.com/vita-epfl/trajnetplusplusdata

- Data+evaluation protocols
- Challenge on Aicrowd
  - https://www.aicrowd.com/challenges/trajnet-a-trajectory-forecasting-challenge

# Current paradigm

Learned Representation

✗ **Not** Robust

Observed sequence
Forecasted sequence by [1]
✗ **Collision**

[1]  Ynet, ICCV'21, Top ranked model in Trajnet++ public challenge

# Current paradigm



Current paradigm:

Use **data** to **learn** the task of **navigation from scratch**

✗ **Not** Robust

Our proposed paradigm:

Use **data** to **only learn concepts we struggle to explain** (*e.g.*, our social norms)

## Because

1. Imbalanced/missing data

## Solution

▪ Knowledge-Data

Alexandre Alahi, EPFL

VITA

# Proposed Knowledge-Data paradigm

Physical laws

Social rules

- - -▶ Output of knowledge

→ Learned residual from data.
Physically constrained space

## Because

1. Imbalanced/missing data

## Solution

- Knowledge-Data
  - Knowledge as input

[1] Injecting knowledge in data-driven vehicle trajectory predictors, **TRC'21**

# Proposed Knowledge-Data paradigm

Alexandre Alahi, EPFL



Physical laws

Social rules

Learned Disentangled Social Representation

## Outcome

✓ Generalizable
(low-shot transfer)

## Because

1. Imbalanced/missing data

## Solution

▪ Knowledge-Data
  • Knowledge as input
  • Knowledge within

VITA

[1] Towards Robust and Adaptive Motion Forecasting: A Causal Representation Perspective, **CVPR'22**

# Proposed Knowledge-Data paradigm

Alexandre Alahi, EPFL



Physical laws

Social rules

Repulsive

Attractive

Learned Disentangled Social Representation

## Because

1. Imbalanced/missing data

## Solution

- Knowledge-Data
  - Knowledge as input
  - Knowledge within
  - Knowledge as supervision

[1] Interpretable Social Anchors, **CVPR'21**

# Knoweldge-data driven mathematical framework

Discrete Choice Models

$$U = V + \varepsilon$$

Utility  Systematic  Random





Predicted choice

[1] Enhancing Discrete Choice Models with Representation Learning, **TRB'20**

# Knoweldge-data driven mathematical framework

Discrete Choice Models    & Representation Learning [1]

$$U \quad = \quad f^{\ attractive} \ + \ f^{\ repulsive} \quad + \quad \varepsilon \quad + \quad r$$



Encoder

Encoder

Encoder

Encoder

Encoder

Encoder

Learned
Representation

[1] Enhancing Discrete Choice Models with Representation Learning, **TRB'20**

# Knoweldge-data driven mathematical framework

Discrete Choice Models    & Representation Learning

$$U_+R = f^{\text{attractive}} + f^{\text{repulsive}} + \varepsilon + \underbrace{r}$$

Learned Representation

Encoder

Encoder

Predicted anchors
R ▶ Neural-network refinement
Final predicted future sequence

[1] Interpretable Social Anchors, **CVPR'21**

# Proposed Knowledge-Data paradigm

Physical laws

Social rules

Repulsive

Attractive

Tasks

Learned Disentangled Social Representation

## Because

1. Imbalanced/missing data

## Solution

- Knowledge-Data
  - Knowledge as input
  - Knowledge within
  - Knowledge as supervision

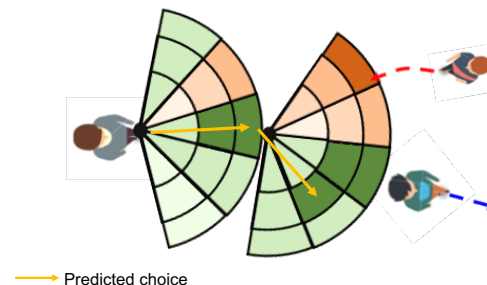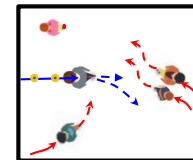# Proposed Knowledge-Data paradigm



Physical laws

Repulsive
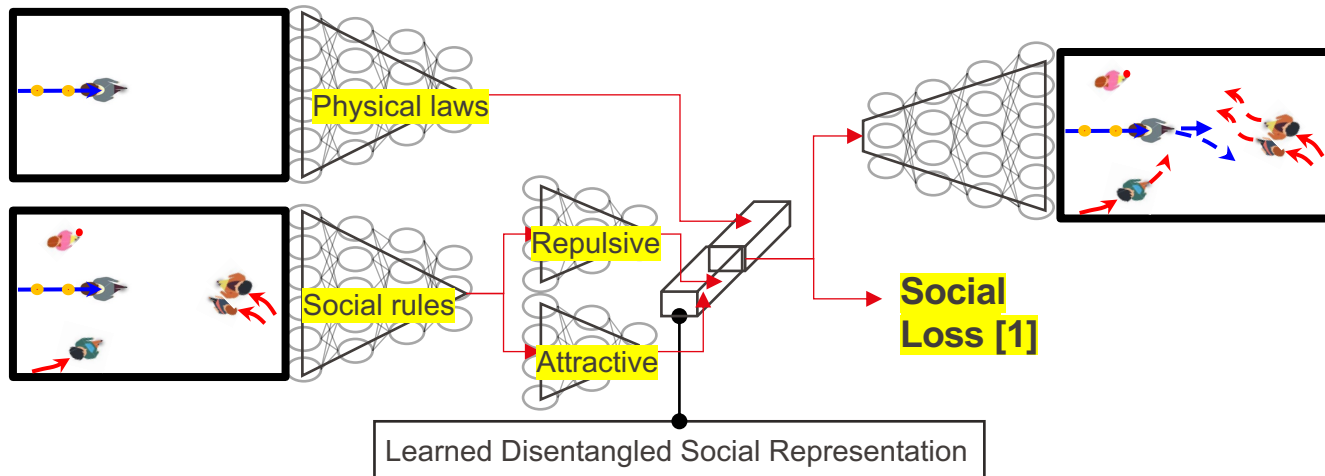
Social rules

Attractive

Social Loss [1]

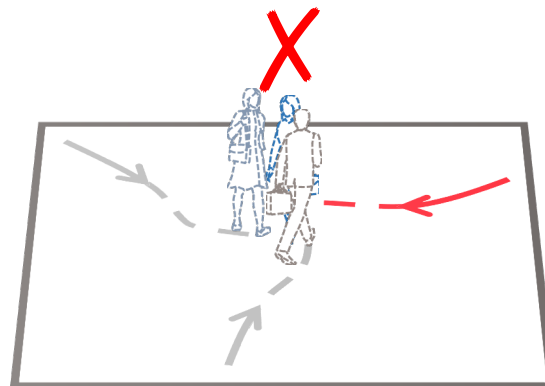Learned Disentangled Social Representation

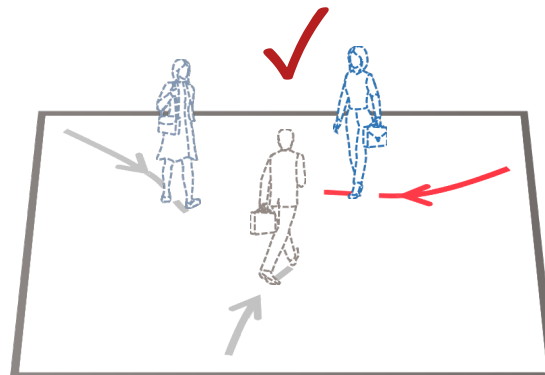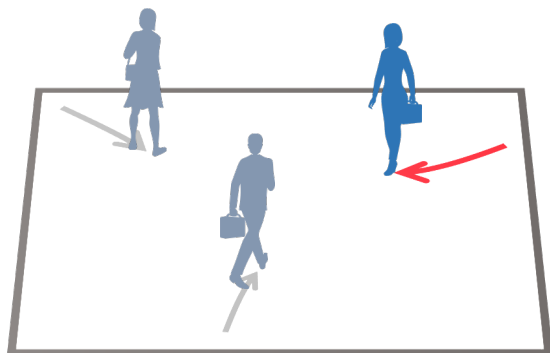## Because

1. Imbalanced/missing data
2. Positive examples only

## Solution

- Knowledge-Data
- w/ Opposite principle

Alexandre Alahi, EPFL

[1] Social NCE, **ICCV'21**

VITA

# Negative data augmentation

History Observation

Future Prediction

[1] Social NCE, **ICCV'21**

# Proposed Knowledge-Data paradigm

**Outcome**

✓ Robust

✓ Generalizable

✓ Interpretable
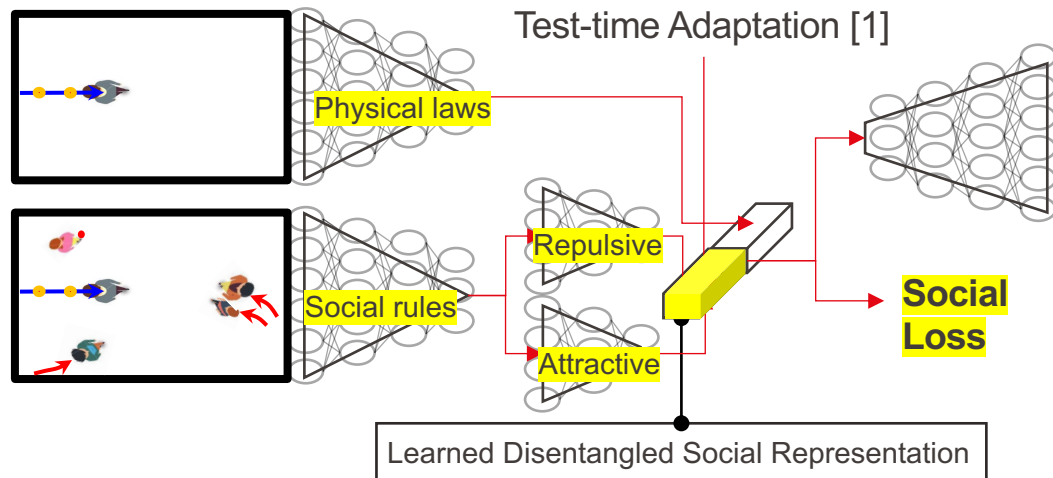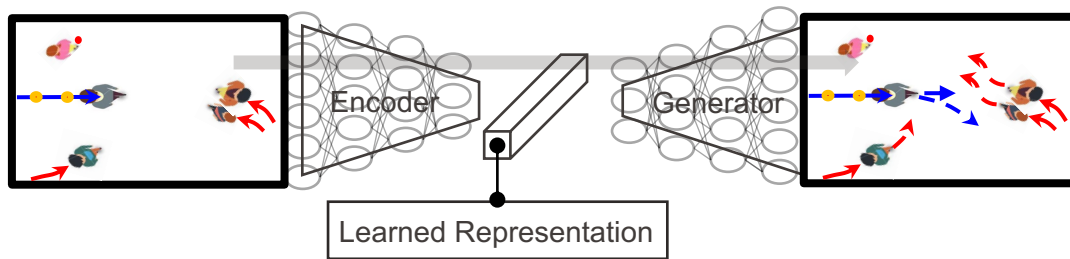
**Because**

1. Imbalanced/missing data
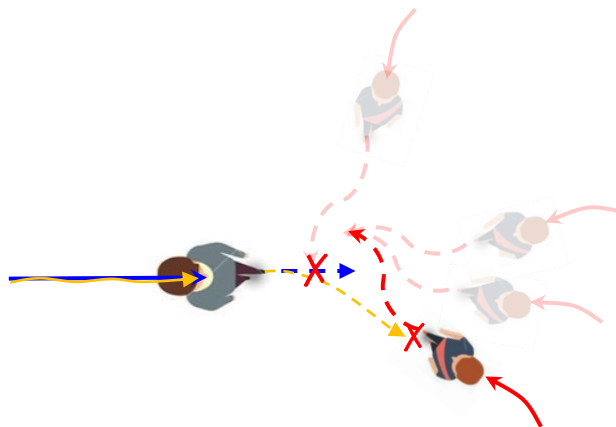2. Positive examples only
3. Distributional shifts

**Solution**

▪ Knowledge-Data

▪ w/ Opposite principle

▪ w/ Low-rank principle

[1] Test Time Training++, **NeurIPS'21**

# New evaluation protocol

Encoder

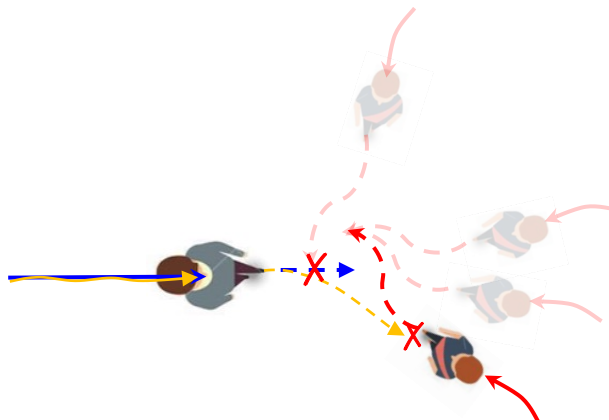Generator

Learned Representation

## Outcome

✓ **New evaluation** based on realistic adverserial examples [1]

✓ **Robust training**



→ Observed sequence
⇢ Forecasted sequence by [2]
→ Perturbed observation by < 7 cm
⇢ Forecasted sequence leading to collision
✗ **Collision**

[1] S-attack library, https://s-attack.github.io/

[2] Ynet, ICCV'21, Top ranked model in Trajnet++ public challenge

# Quantitative results

| Baseline | Original collision rate |
|---|---|
| S-LSTM (CVPR'16) | 7.8% |
| S-Att (ICRA'18) | 9.4% |
| S-GAN (CVPR'18) | 13.9% |
| D-Pool (ITS'2021) | 7.3% |
| S-STGCNN (CVPR'20) | 16.3% |
| PECNet (ECCV'20) | 15.0% |

=> 6.5% w/ aug

VITA

[1] S-attack library, https://s-attack.github.io/

# Qualitative results

D-Pool[ITS'2021]

S-STGCNN[CVPR'20]

S-GAN[CVPR'18]

[1] S-attack library, https://s-attack.github.io/

# New evaluation protocol

## Outcome

✓ **New evaluation** based on realistic adverserial examples [1]

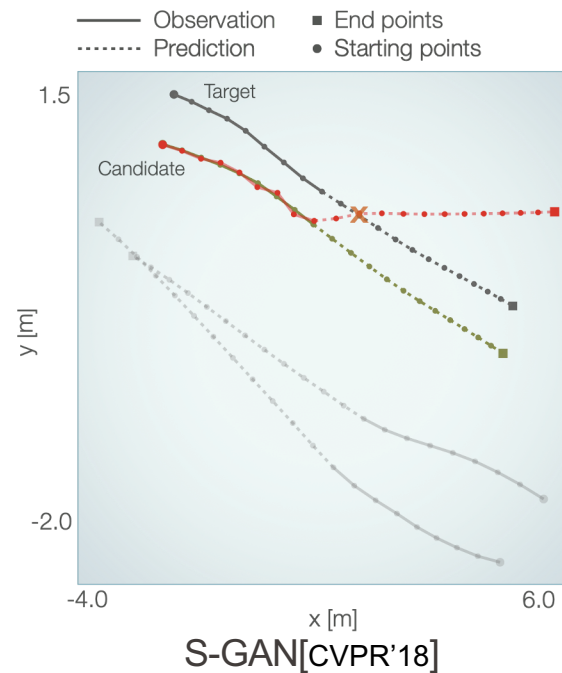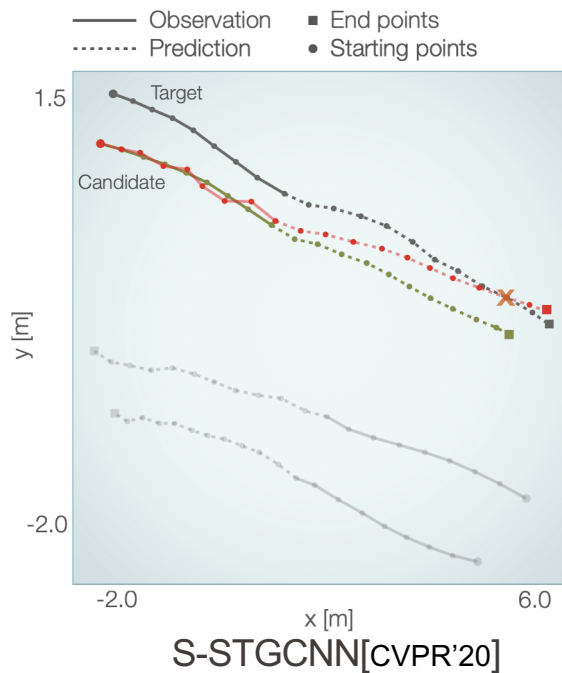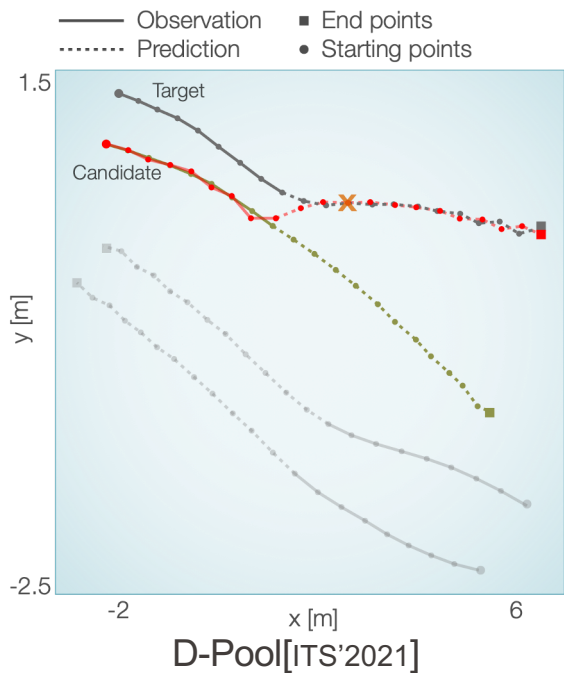✓ **Robust training**



Encoder

Generator

Learned Representation

→ Observed sequence
⇢ Forecasted sequence by [2]
→ Perturbed observation by < 7 cm
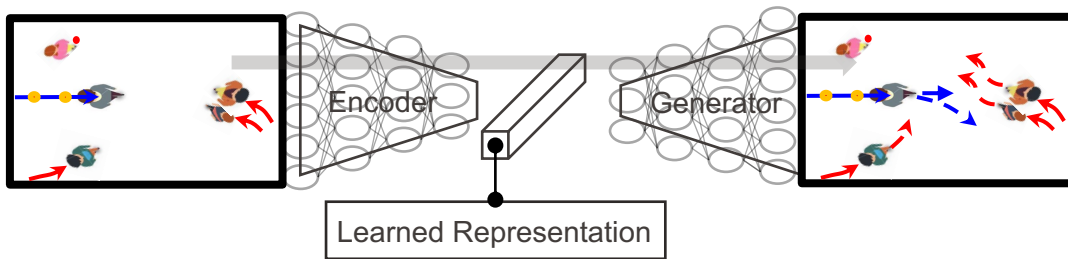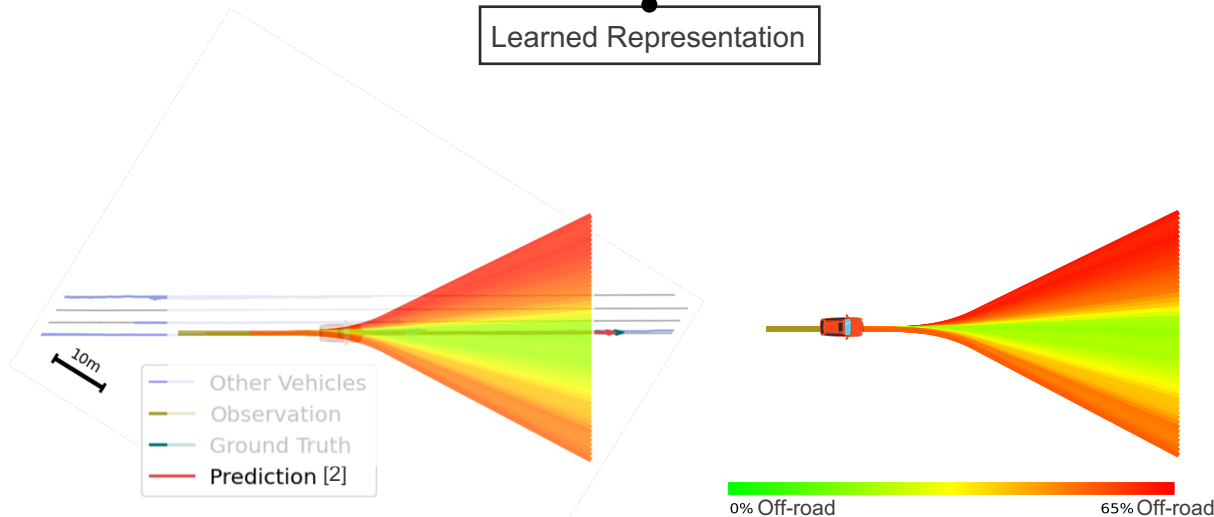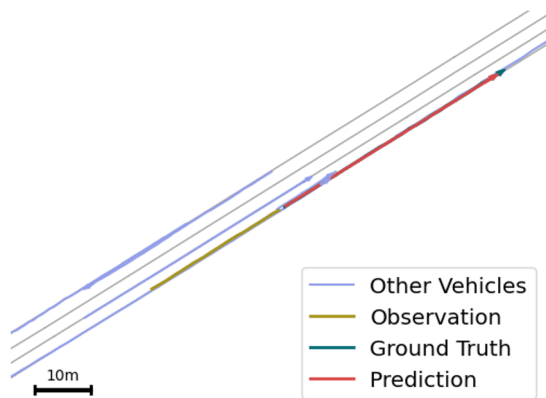⇢ Forecasted sequence leading to collision
✗ **Collision**

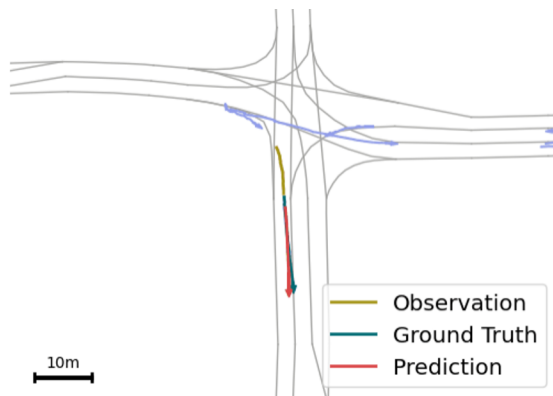[1] S-attack library, https://s-attack.github.io/

[2] Ynet, ICCV'21, Top ranked model in Trajnet++ public challenge

# New evaluation protocol

Alexandre Alahi, EPFL



Encoder

Generator

Learned Representation

## Outcome

✓ **New evaluation** based on realistic adverserial examples [1]

✓ **Robust training**



10m

Other Vehicles
Observation
Ground Truth
**Prediction** [2]

0% Off-road                    65% Off-road



VITA

[1] Vehicle trajectory prediction works, but not everywhere, **CVPR'22**      [2] LaneGCN, ECCV'20, Top ranked model in Argoverse public challenge

# Scene generation

✓ Atomic scene generation functions



| Simple turn | Double turn | Ripple road |

[1] Vehicle trajectory prediction works, but not everywhere, **CVPR'22**

# Quantitative results

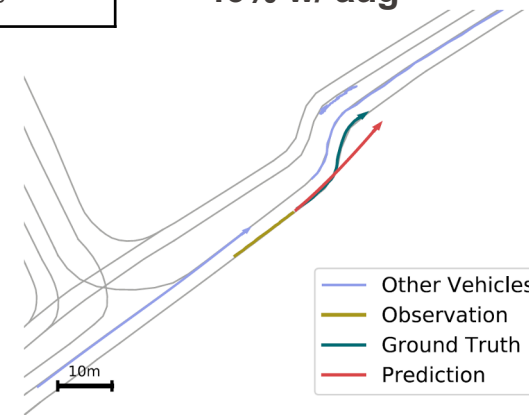| Baseline | Original off-road | Generated (ours) off-road |
|---|---|---|
| DATF (ECCV20) | 2% | 82% |
| WIMP (arXiv20) | 1% | 63% |
| LaneGCN (ECCV'20) | 1% | 66% |

=> 46% w/ aug



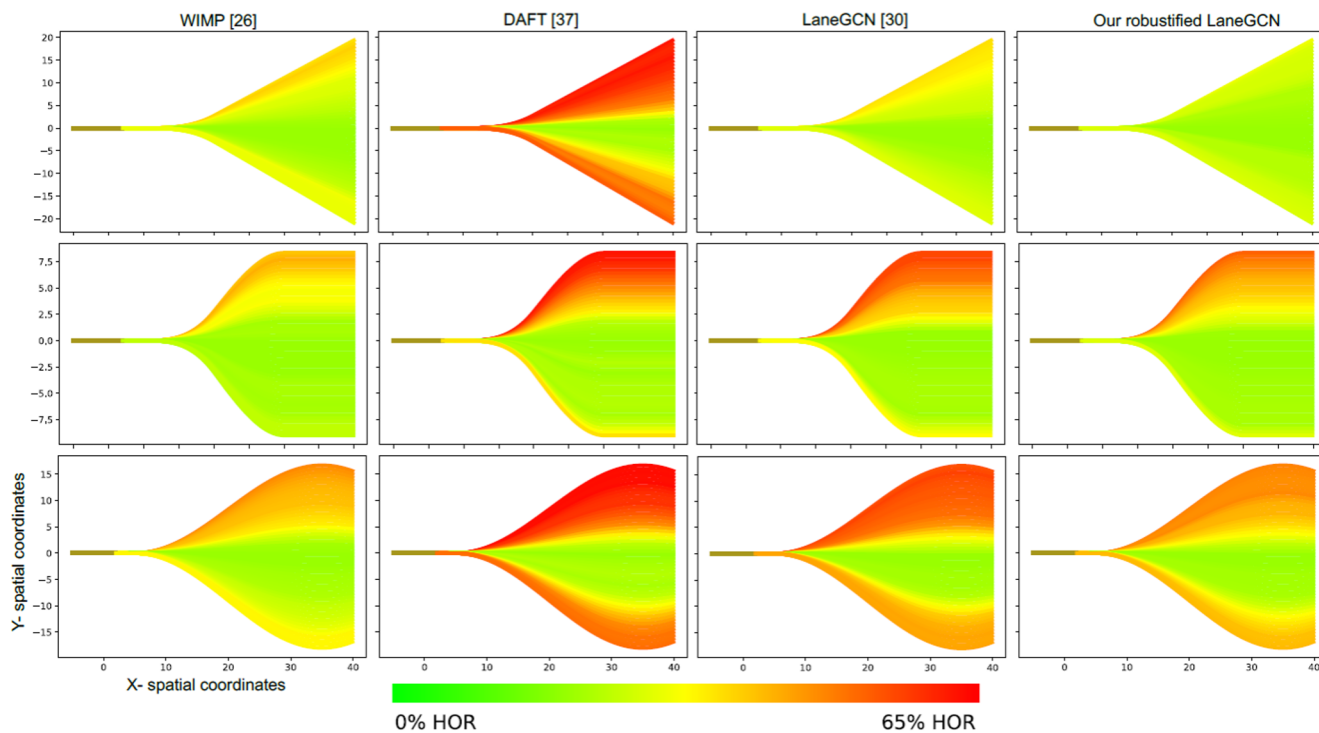(a) DATF     (b) WIMP     (c) LaneGCN
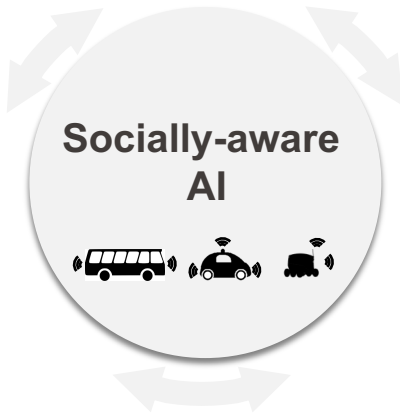
[1] Vehicle trajectory prediction works, but not everywhere, **CVPR'22**

# Discussions

[1] Vehicle trajectory prediction works, but not everywhere, **CVPR'22**

Social Forecasting

Perceiving

**Socially-aware AI**

Planning

# Crowd-Robot Interaction [1]

Alexandre Alahi, EPFL



Human-Robot Interaction [1,2]
&
Our Human-Human Interaction [3]

Previous works
[1] HRI, Chen, C., *et al.,*          IROS'17
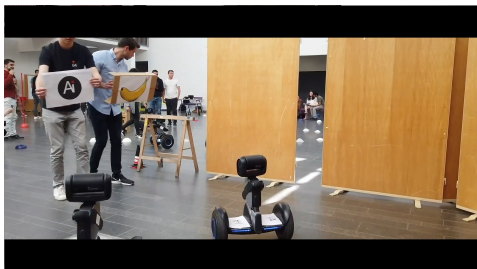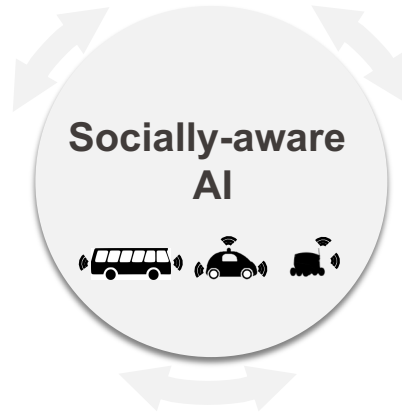[2] HRI, Everett, M., *et al.,*          IROS'18
Our work
[3] Crowd-Robot Interaction,          ICRA'19

VITA

# Human-Robot Tandem Race

Social Forecasting

Perceiving

**Socially-aware AI**

Planning

Thank you!

# #Open Science

Alexandre Alahi, EPFL

**Perception:**
[1] S. Kreiss et al., OpenPifPaf **library** for pose estimation, **CVPR'19, ICCV'21 (licensed)**
[2] L. Bertoni et al., 3D perception **library**, **ICCV'19, ICRA'21**
[3] L. Bertoni et al., Perceiving Social Distancing, **ITS'20**
[4] G. Adaimi et al., Deep Visual Re-identification with Confidence, **TRC'21**
[5] T. Mordan et al., Detecting 32 human attributes, **ITS'21**
**Prediction:**
[6] Kothari et al., Trajnet++ **library** for spatio-temporal forecasting tasks (>15 implemented models)
[7] Kothari et al., Social Anchor, **ICCV'21**
[8] Liu et al., Social NCE, **ICCV'21**
**Planning:**
[9] C. Chen et al., Crowd-Robot Interaction, **ICRA'19**
**Generative models:**
[10] Y. Liu* et al., Collaborative Sampling in GAN, **AAAI'20**
[11] A. Carlier et al., Deep SVG, **NeurIPS'20**
**DCM + NN**
[12] B. Sifringer et al., L-MNL, **TRB'20**
**Test-time training:**
[13] Y. Liu* et al., TTT++, **NeurIPS'21**
**Tools**
[14] Video Ultimate labeling
[15] S-attack library**, CVPR'22**

VITA

**GitHub**

Code on-line: **vita.epfl.ch/code**