# Perceiving and Predicting Human behaviours in the built environments

Prof. Alexandre Alahi

- Assistant Professor at EPFL since 2017
- Director of the VITA lab
- 5 years at Stanford University
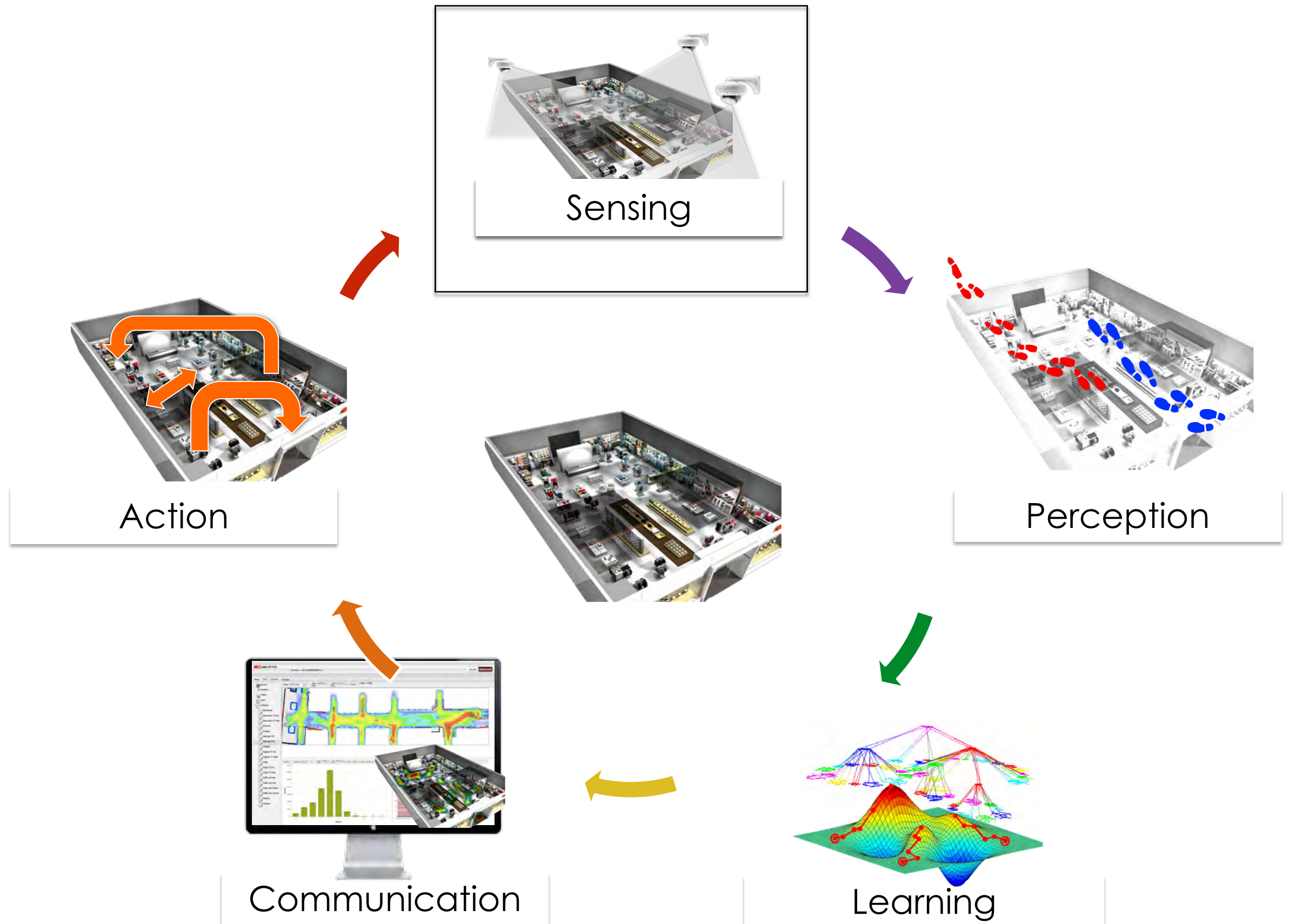- Founded and advise startups (in retails)

**VITA**

#Open Science

# "Man is by nature a social animal" - Aristotle

How can machines learn **human-human** interactions?

How can machines learn **human-space** interactions?

# AI for the built environments



Sensing

Perception

Learning

Communication

Action

2 Terminals

132 sensors

>20,000 m$^2$

24/7 over a year

>42 million humans

(Lausanne, Basel)

# A corridor with 32 sensors

## Top View

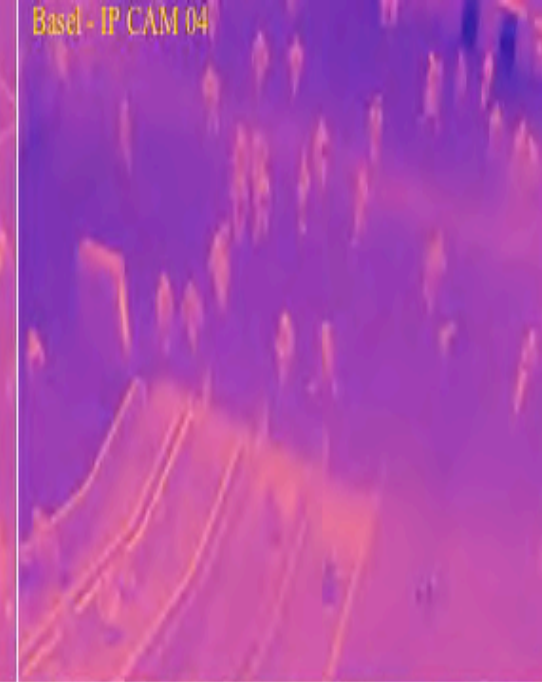# Collecting long-term trajectories

Our data (>42 million humans)

# AI for the built environments



Sensing

Perception

Action

Communication

Learning

# A corridor with 32 sensors

## Top View



10 m
33 ft. / 100 m / 330 ft.

Alex Alahi –ENAC/IIC/VITA

# Perceiving

# Perceiving
# Socially-aware cues



## 3D Body poses + Activities + Relationships

Walking

Standing Standing

Walking

Talking

2.3m

2.1m

2.2m

4.5m

Our work
[1] PifPaf: Composite Fields for Human Pose Estimation, **CVPR'19** (Live Demo: https://vitademo.epfl.ch)
[2] Convolutional Relational Machine for Group Activity Recognition, **CVPR'19**
[3] Monoloco: Monocular 3D pedestrian localization and uncertainty estimation, **ICCV'19**

12

# Computer vision for the built environments

3 Challenges :
1- Limited resolution with partial information
2- Efficiency (Real-time)
3- Many tasks with unbalanced labels



Encoder

Z

Generator

Perceive current state

Body poses    + Activities

Walking    Talking

Shared representation for

Perception

# Computer vision for the built environments

3 Challenges :
1- Limited resolution with partial information
2- Efficiency (Real-time)
3- Many tasks with unbalanced labels



Encoder

$z$

Generator

Perceive current state

Body poses + Activities

Walking Talking

Shared backbone

$\frac{d\mathcal{L}_\mathbf{S}}{dz}$

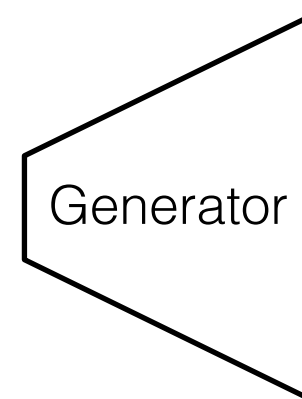Task 1 → $L_{Task\ 1}$

Task 2 → $L_{Task\ 2}$

Task 3 → $L_{Task\ 2}$

$\frac{d\mathcal{L}_{\mathbf{F}_A}}{dz}$

Task n → $L_{Task\ n}$

Where L=loss

Composite fields formalism

A Feature map can jointly encode scalars and vectors

# Jointly Perceiving 32 Attributes on Pedestrian's Appearance, Behaviour, Intention



[1] Detecting 32 Pedestrian Attributes for Autonomous Vehicles, arxiv

# Perceiving Intentions

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



[1] Pedestrian Intention Prediction: A Convolutional Bottom-Up Approach, arxiv

# Object
## Detection & Tracking
2 steps

# Semantic keypoints
## Joint Detection & Tracking



[1] PifPaf: Composite Fields for Human Pose Estimation, **CVPR'19**

(Live Demo: https://vitademo.epfl.ch)   19

# Perceiving
# 3D body joints

3D Body joints



[1] L. Bertoni *et al.*, MonoLoco, ICCV'19

[2] W. Deng et al., Joint Human Pose Estimation and Stereo Localization, ICRA'20

# AI for the built environments



Sensing

Perception

Learning

Communication

Action

# Social Forecasting

Input:      Given a sequence of states, *e.g.*, $(x^t, y^t)$ coordinates in time

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



**Challenge 1:** Social Interactions

**Challenge 2:** free will

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

→ Observed sequence

⇢ Forecasted sequence

# Social Forecasting

- Previous works

| Knowledge-driven |
| :---: |

- Social Forces Model [1],

$$F = F^{\text{attractive}} + F^{\text{repulsive}} \ldots$$



- Discrete Choice Model [2]

$$\underbrace{U}_{\text{Utility}} = \underbrace{V}_{\text{Systematic}} + \underbrace{\varepsilon}_{\text{Random}}$$

✓ Interpretability
X Predictability

Previous works
[1] Helbing *et al.*, Physical review, '95
[2] Antonini *et al.*, Transportation Research, '06

# Social Forecasting

- Previous works

- Our works

| Knowledge-driven |
|---|

| Data-driven methods |
|---|

- Social Forces Model [1],

$$F = F^{\text{attractive}} + F^{\text{repulsive}} \dots$$

- Discrete Choice Model [2]

$$\underbrace{U}_{\text{Utility}} = \underbrace{V}_{\text{Systematic}} + \underbrace{\varepsilon}_{\text{Random}}$$

Encoder    Generator

✓ Interpretability
X Predictability

X Interpretability
✓ Predictability

Previous works
[1] Helbing *et al.*, Physical review, '95
[2] Antonini *et al.*, Transportation Research, '06
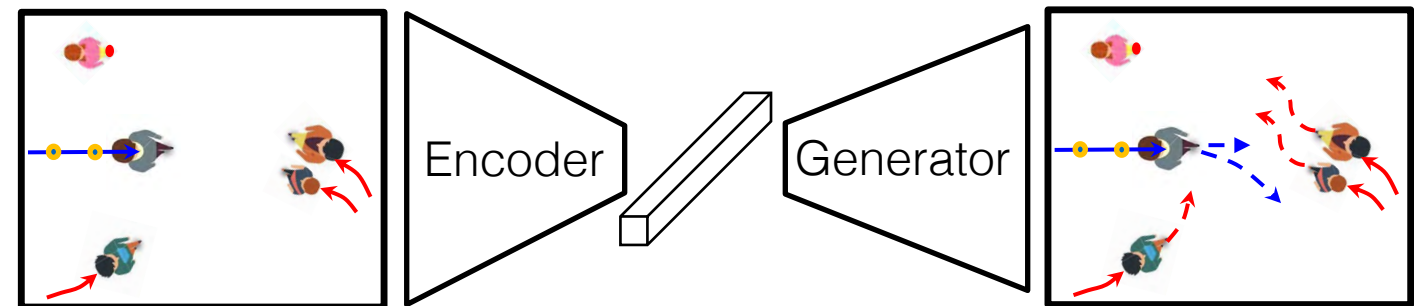
# Social Forecasting
## Learning Representation of Crowds

Challenge 2: Free will



GAN [1]



Average L2 Loss

$$\ell_{Recon} = \frac{1}{K} \sum_k \| Y_i - \hat{Y}_i^{(k)} \|$$

Winner-takes-all Loss

$$\mathcal{L}_{variety} = \min_k \| Y_i - \hat{Y}_i^{(k)} \|$$

Adversarial Loss [1]

Adversarial Loss w/ Collab. Sampling [2]

**Our work**
[1] Social GAN (Generative Adversarial Network),      CVPR'18
[2] Collaborative Sampling in Generative Adversarial Networks, AAAI'20

26

# Performance evaluation

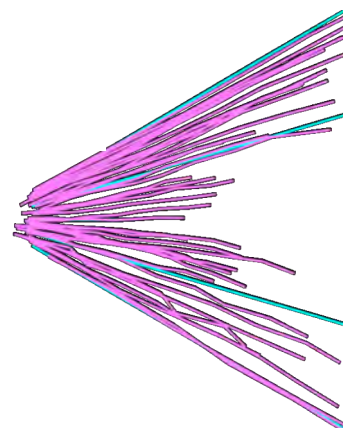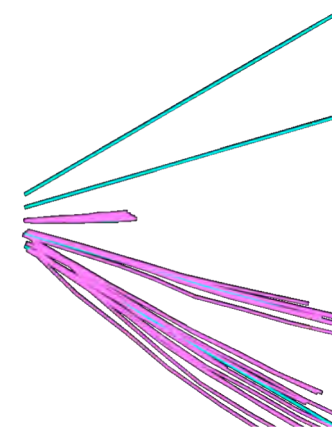| | Methods* | ADE/FDE | |
|---|---|---|---|
| Knowledge-driven | Kalman Filter | 0.87/1.69 | |
| | DCM '06 | 0.67/1.43 | |
| | Social Force '98 | 0.68/1.40 | 0 |
| | ORCA '08 | 0.68/1.40 | |
| Our data-driven | LSTM '14 | 0.61/1.31 | |
| | S-GAN '18 | 0.57/1.24 | |
| | S-LSTM '16 | 0.54/1.17 | 7 |
| | D-LSTM '20 | 0.57/1.24 | |

0.55/1.18    7.6

ADE: Average displacement error in m
FDE: Final displacement error in m
All references available in [1]

Our work
[1] Human trajectory forecasting: a deep learning perspective, arxiv

# AI for the built environments



Sensing

Perception

Learning

Communication

Action

# What can we learn from all these trajectories?



| # People | Av. duration | Av distance | Density (up to) | # Paths (O/D) |
|----------|--------------|-------------|-----------------|---------------|
| 42 million | 1 minute | 100m | 1 pedestrian/m$^2$ | 196 |

Bertoni, L., Kreiss, S., Alahi,A,
Perceiving humans: from monocular 3D localization to social
distancing,
arxiv

Thank you

VITA

#Open Science

# #Open Science

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

**Perception:**
[1] S. Kreiss et al., OpenPifPaf **library** for pose estimation, **CVPR'19 (licensed)**
[2] L. Bertoni et al., Monocular 3D Pedestrian Localization and Uncertainty Estimation, **ICCV'19**
[3] L. Bertoni et al., MonStereo, Stereo 3D detection
[4] L. Bertoni et al., Perceiving Social Distancing, **ITS'20**
[5] G. Adaimi et al., Perceiving Traffic from Aerial Images
[6] G. Adaimi et al., Deep Visual Re-identification with Confidence

**Prediction:**
[7] Kothari et al., Trajnet++ **library** for spatio-temporal forecasting tasks (>15 implemented models)

**Planning:**
[8] C. Chen et al., Crowd-Robot Interaction: Crowd-aware Robot Navigation with Attention-based Deep Reinforcement Learning, **ICRA'19**

**Generative models:**
[9] Y. Liu* et al., Collaborative Sampling in GAN, **AAAI'20**
[10] A. Carlier et al., Deep SVG, **NeurIPS'20**

**DCM + NN**
[11] B. Sifringer et al., L-MNL, **TRB'20**

**Tools**
[12] Video Ultimate labeling