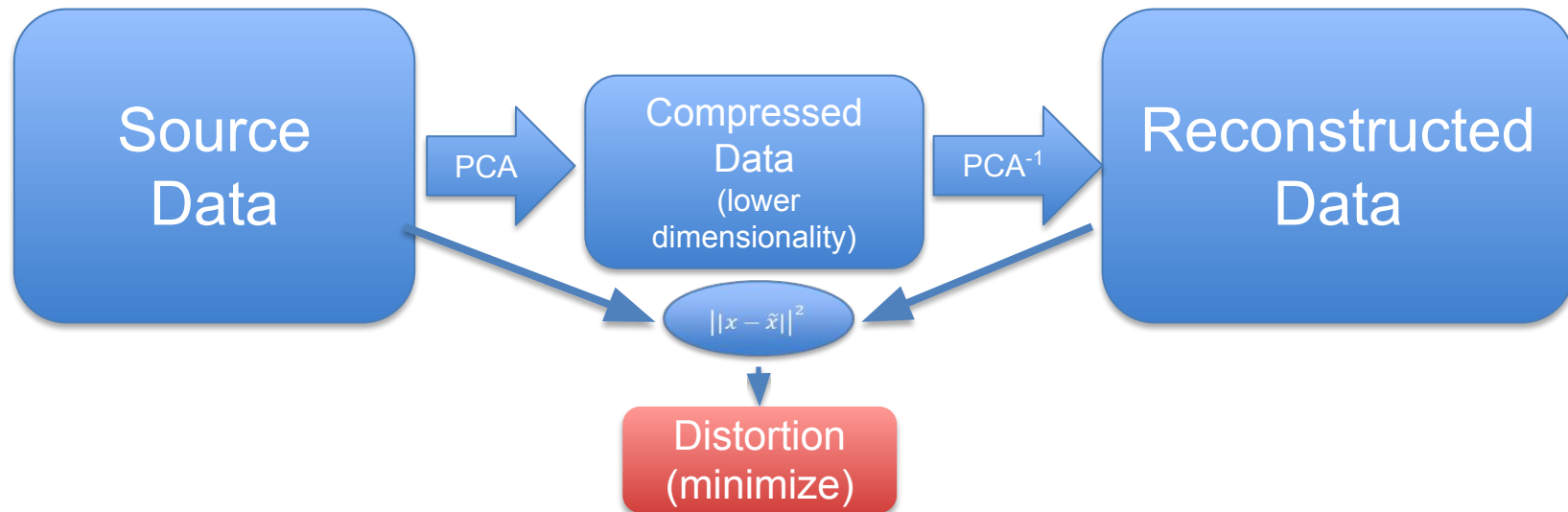


Principal Component Analysis

Machine Learning 2022-23
UML Book Chapter 23

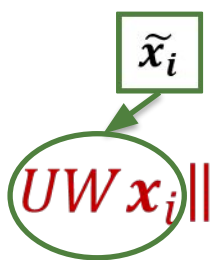
Dimensionality Reduction



- Take data from an **highly** dimensional space and project to a **lower** dimensional one
- Many applications: reduce number of features (learn with less samples, lower computation req.), capture most important aspects of the data for subsequent analysis, visualization, etc..
- Lower dimensional data should be a **good approximation** of the higher dimensional representations
- Good approximation**: minimize error obtained by reprojecting the data back to the high dimensional space → similar to lossy data compression
- Focus on **linear mapping** of the data (represented by a matrix multiplication)
- Principal Component Analysis (PCA)**: find the linear mapping that minimizes the mean squared error in the reprojection

Principal Component Analysis (PCA)

- ❑ $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathbb{R}^d$: data points
- ❑ $W \in \mathbb{R}^{n,d}$ ($n < d$): mapping $\mathbf{x} \rightarrow \mathbf{y} = W\mathbf{x}$
 - where $\mathbf{y} = W\mathbf{x} \in \mathbb{R}^n$ is a lower dimensional representation of $\mathbf{x} \in \mathbb{R}^d$
- ❑ $U \in \mathbb{R}^{d,n}$ ($n < d$): inverse mapping $\mathbf{y} \rightarrow \tilde{\mathbf{x}} = U\mathbf{y}$
 - used to recover an approximation $\tilde{\mathbf{x}} = UW\mathbf{x}$ of \mathbf{x}
- ❑ Target: find the lower dimensional representation that better approximates the data \rightarrow that leads to minimum squared distance between $\tilde{\mathbf{x}}$ and \mathbf{x}

$$\operatorname{argmin}_{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}} \sum_{i=1}^m \|\mathbf{x}_i - UW\mathbf{x}_i\|_2^2$$


PCA : Algorithm Idea

Target:

seek for the n -dimensional basis that best captures the variance in the d -dimensional data

Procedure:

1. First principal component ($p.c.$) = direction with largest projected variance
 2. Second $p.c.$ = orthogonal direction with largest projected variance
 - i.e., largest remaining variance after removing the first $p.c.$
 3. ...(iterate for $3 \dots n$)....
- First derive the 1-dimensional subspace that maximizes the projected variance and then proceed iteratively

Lemma

There exist an optimal solution (U^*, W^*) of $\operatorname{argmin}_{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}} \sum_{i=1}^m \|x_i - UWx_i\|_2^2$ where:

Recall: orthonormal

- the columns of U^* are orthonormal (i.e., $(U^*)^T U^* = I$)
- $W^* = (U^*)^T$

- $u_i^T u_j = 0, \forall i \neq j$
- $\|u_i\| = 1 = u_i^T u_i, \forall i$

Demonstration:

1. Fix U, W and consider the mapping $x \rightarrow UWx$
 - The range of the mapping is $R = \{UWx ; x \in \mathbb{R}^d\}$
2. $V \in \mathbb{R}^{d,n}$: matrix whose column form an orthonormal basis of R
 - Recall that $V^T V = I$ and $\forall x \in R: \exists y \in \mathbb{R}^n$ with $x = Vy$
3. $\forall x \in \mathbb{R}^d, \forall y \in \mathbb{R}^n$:
 - $\|x - Vy\|_2^2 = \|x\|^2 + y^T V^T V y - 2y^T V^T x = \|x\|^2 + \|y\|^2 - 2y^T (V^T x)$
4. Minimize $\|x\|^2 + \|y\|^2 - 2y^T (V^T x)$ w.r.t y : set $\nabla = 0 \rightarrow 2y - 2(V^T x) = 0 \rightarrow y_{opt} = V^T x$
5. $\forall x: \operatorname{argmin}_{\tilde{x} \in R} \|x - \tilde{x}\|_2^2 = Vy_{opt} = V(V^T x)$: it is the best approximation in subspace R
6. $\forall x$: includes also x_1, \dots, x_m (data vectors): $\sum_{i=1}^m \|x_i - UWx_i\|^2 \geq \sum_{i=1}^m \|x_i - VV^T x_i\|^2$, so we can replace U, W with VV^T without increasing the objective
7. Holds for $\forall U, W$: there exist a solution that minimize $\sum_{i=1}^m \|x_i - UWx_i\|^2$ with V orthonormal columns and $W = U^T$

Optimization Problem

There exist an optimal solution (U^*, W^*) of $\operatorname{argmin}_{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}} \sum_{i=1}^m \|x_i - UWx_i\|_2^2$ where:

- the columns of U^* are orthonormal (i.e., $(U^*)^T U^* = I$)
- $W^* = (U^*)^T$

The optimization problem can be rewritten as:

$$\operatorname{argmin}_{U \in \mathbb{R}^{d,n}: U^T U = I} \sum_{i=1}^m \|x_i - UU^T x_i\|_2^2$$

- Trace: Σ elements on diagonal
- It is a scalar
- $\operatorname{trace}(A^T B) = \operatorname{trace}(AB^T) = \operatorname{trace}(B^T A) = \operatorname{trace}(BA^T)$

With some manipulations: $\|x - UU^T x\|_2^2 = \|x\|^2 - 2x^T UU^T x + x^T UU^T UU^T x = \|x\|^2 - x^T UU^T x = \|x\|^2 - \operatorname{trace}(xUU^T x) = \|x\|^2 - \operatorname{trace}(U^T x x^T U)$

$$\operatorname{argmax}_{U \in \mathbb{R}^{d,n}: U^T U = I} \operatorname{trace} \left(U^T \sum_{i=1}^m x_i x_i^T U \right) = \operatorname{argmax}_{U \in \mathbb{R}^{d,n}: U^T U = I} \operatorname{trace}(U^T A U)$$

Notice: $A = \sum_{i=1}^m x_i x_i^T$ is symmetric and positive semidefinite. It can be rewritten as $A = VDV^T$ where D is diagonal (with eigenvalues $D_{d,d} \geq 0$) and $V^T V = VV^T = I$ (the columns of V are the eigenvectors of A)

Theorem (PCA)

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ be arbitrary vectors in \mathbb{R}^d

let $A = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T$

let $\mathbf{u}_1, \dots, \mathbf{u}_n$ be n eigenvectors of A corresponding to the largest n eigenvalues of A



Then a solution of the PCA optimization $\underset{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}}{\operatorname{argmin}} \sum_{i=1}^m \|\mathbf{x}_i - UW\mathbf{x}_i\|_2^2 = \underset{U \in \mathbb{R}^{d,n}: U^T U = I}{\operatorname{argmax}} \operatorname{trace}(U^T A U)$
is to set U to be the matrix whose columns are $\mathbf{u}_1, \dots, \mathbf{u}_n$ and to set $W = U^T$

Notes:

- ❑ Recall: Decompose A as VDV^T (SVD decomposition, D diag. and $V^T V = VV^T = I$)
- ❑ It is a common practice to "center" the examples before applying PCA (i.e., subtract the mean)
- ❑ Computation time is $O(d^3) + O(md^2)$ (the first term for calculating eigenvalues and the second for constructing A)
- ❑ Trick for faster solution in case $d \gg m$ (not part of the course)

Pseudocode

PCA

input

A matrix of m examples $X \in \mathbb{R}^{m,d}$

number of components n

if ($m > d$)

$$A = X^T X$$

Let $\mathbf{u}_1, \dots, \mathbf{u}_n$ be the eigenvectors of A with largest eigenvalues

else

$$B = X X^T$$

Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ be the eigenvectors of B with largest eigenvalues

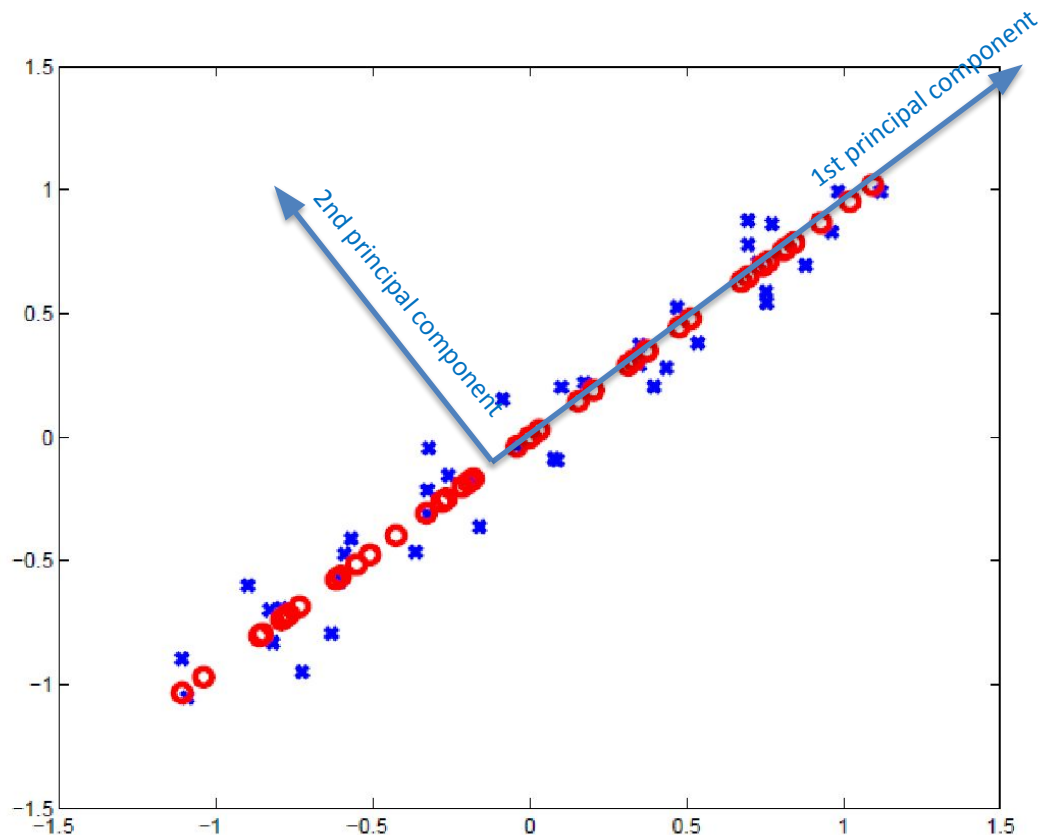
for $i = 1, \dots, n$ set $\mathbf{u}_i = \frac{1}{\|X^T \mathbf{v}_i\|} X^T \mathbf{v}_i$

output: $\mathbf{u}_1, \dots, \mathbf{u}_n$

grayed part: trick for
 $d \gg m$

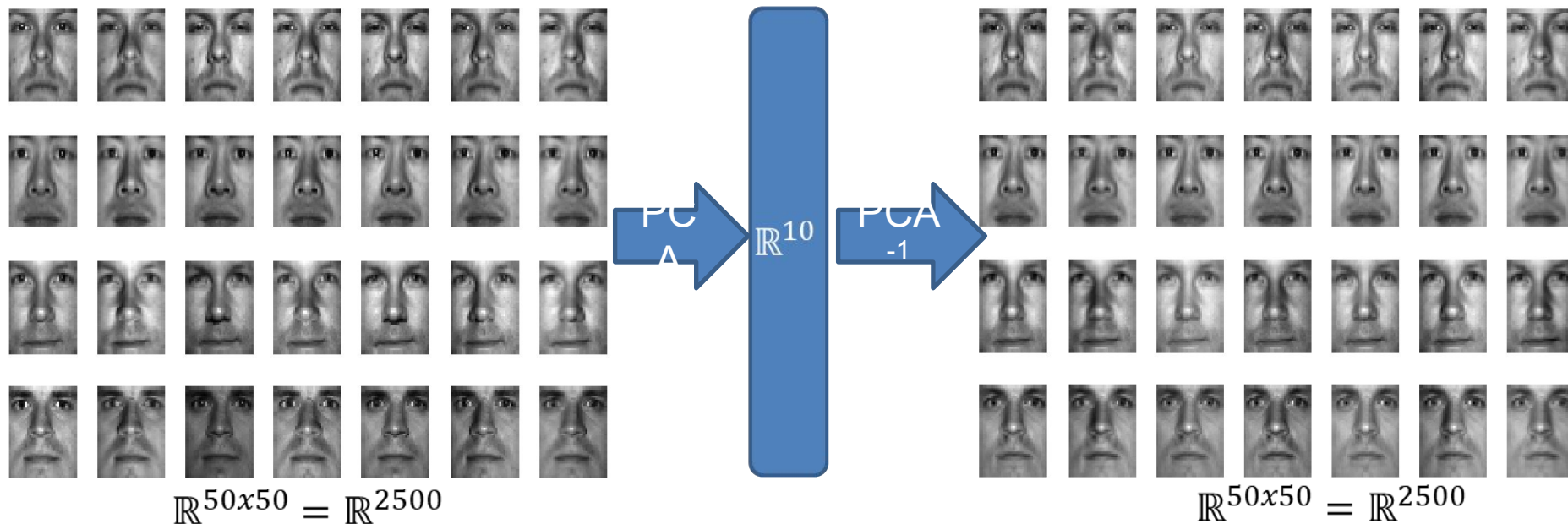
(not part of the course)

Example: From 2D to 1D



*Set of 2D vectors (**blue**) and their reconstruction (**red**)
after dimensionality reduction to 1D with PCA*

Example: Face Compression



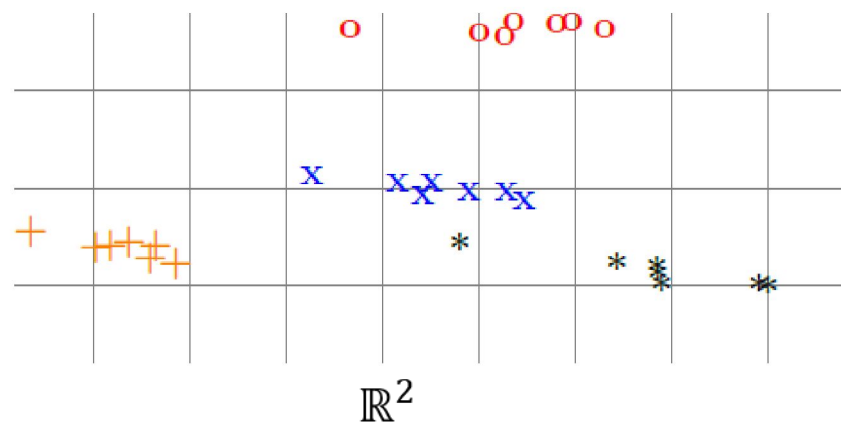
*enlarged
example*



Example: Face Recognition



$$\mathbb{R}^{50 \times 50} = \mathbb{R}^{2500}$$



- Faces with the same type of mark belong to the same individual
- PCA can be used for face recognition ! (*eigenfaces* algorithm)

Logistic info: January

- ❑ Wednesday, January 11: exercises + sample exam solution
- ❑ Friday, January 13: Lab 4 (Keras tutorial) - more information coming soon on Elearning
- ❑ Friday, January 20: exercises + Q&A

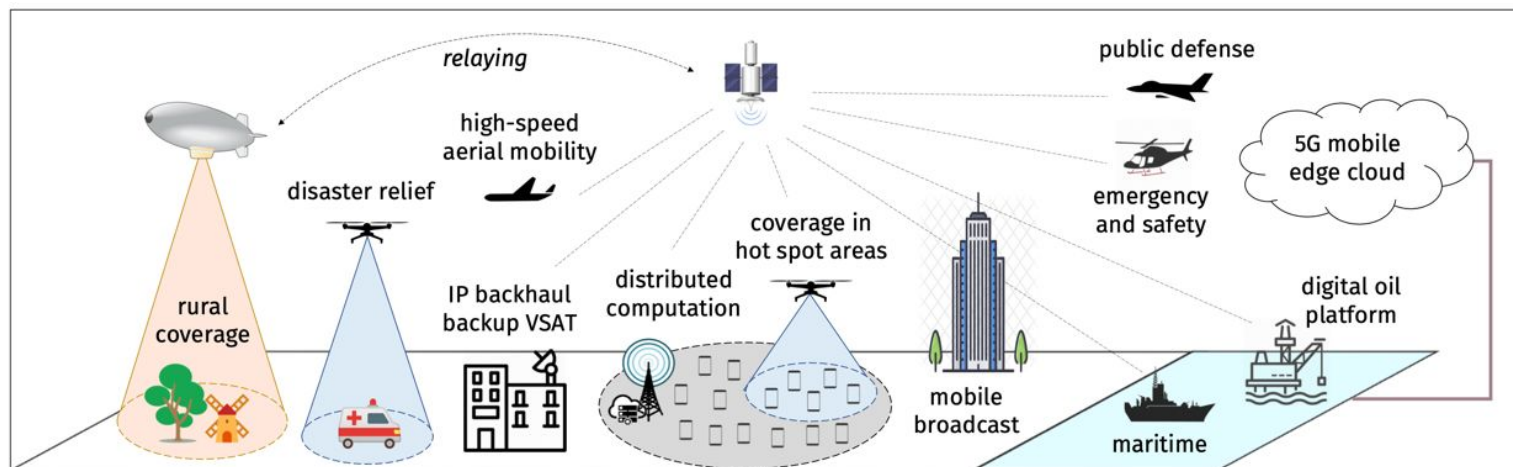
You can always ask for help or a meeting, but response times might be long during the holidays

SIGNET thesis projects

- ❑ Could be useful for thesis or course projects
 - The scope of the project can be tuned
- ❑ You can use ML in a real research context
- ❑ You can also propose your own ideas!

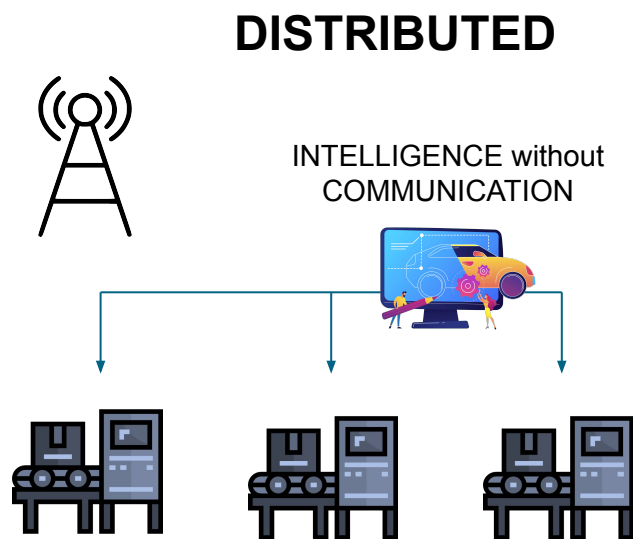
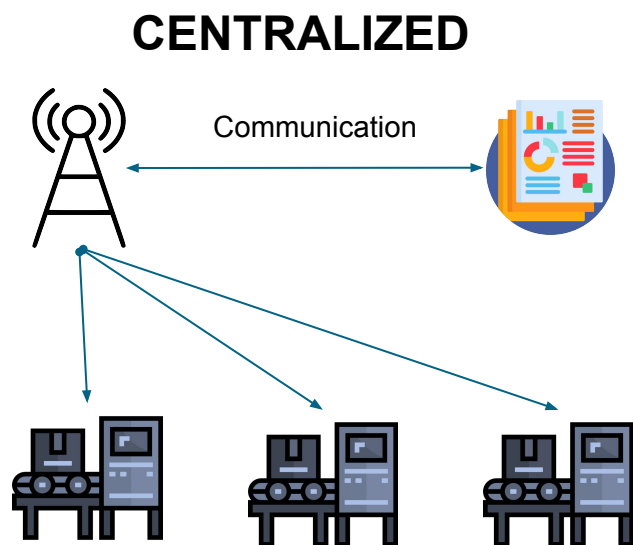
Non-Terrestrial Networks

- ❑ Complementing terrestrial infrastructures with aerial nodes (drones, satellites, high altitude platforms, etc.)
- ❑ ML optimization: how do we distribute comm/computation tasks?



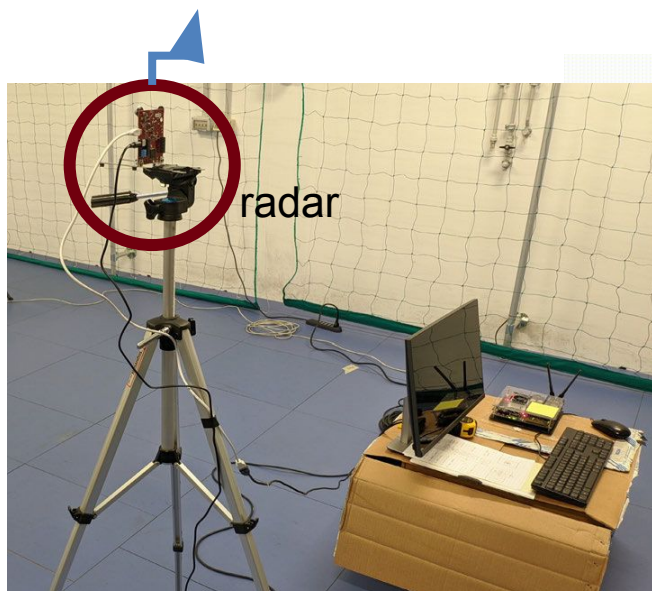
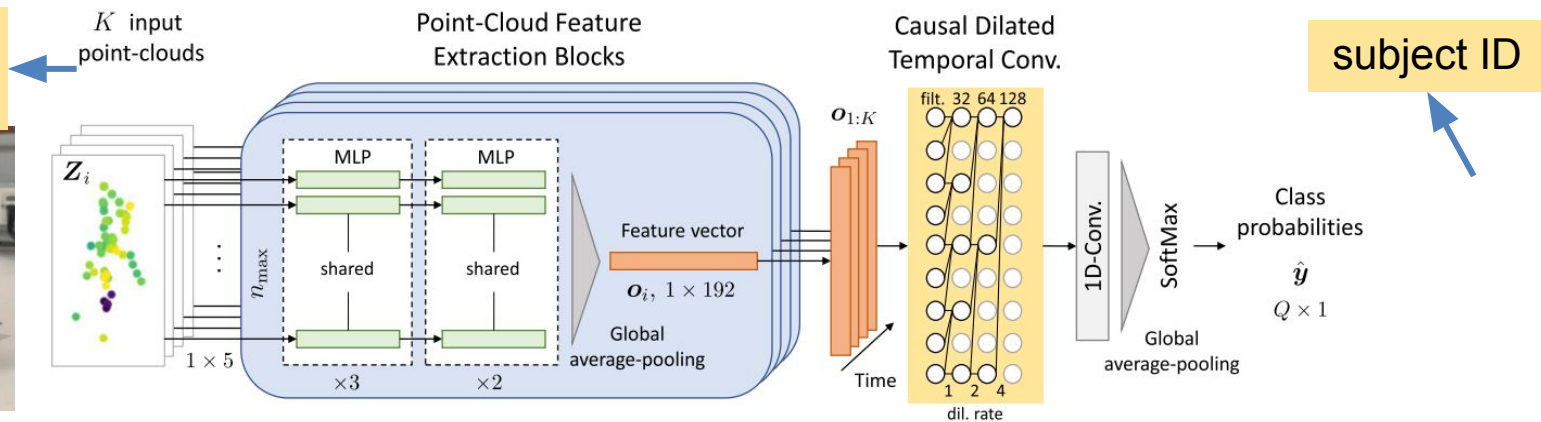
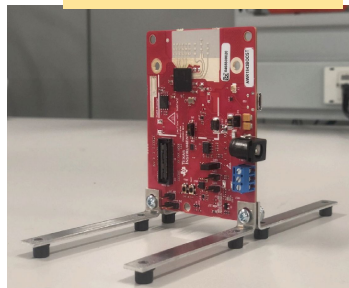
URLLC: extremely low latency services

How do we distribute intelligence to meet the deadlines?

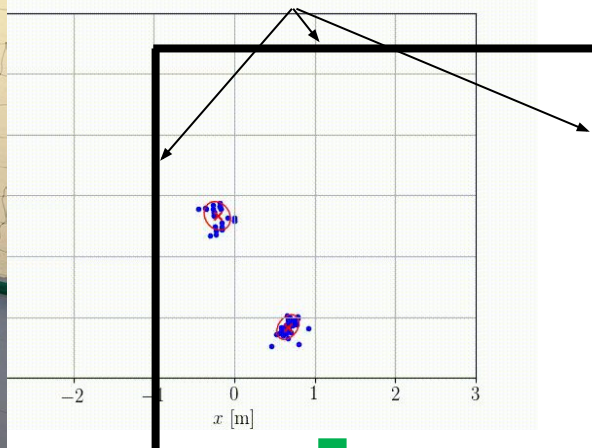


Radar identification

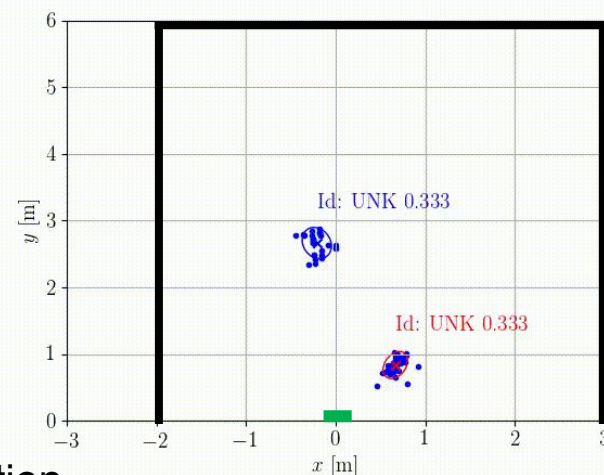
radar raw data



Room walls

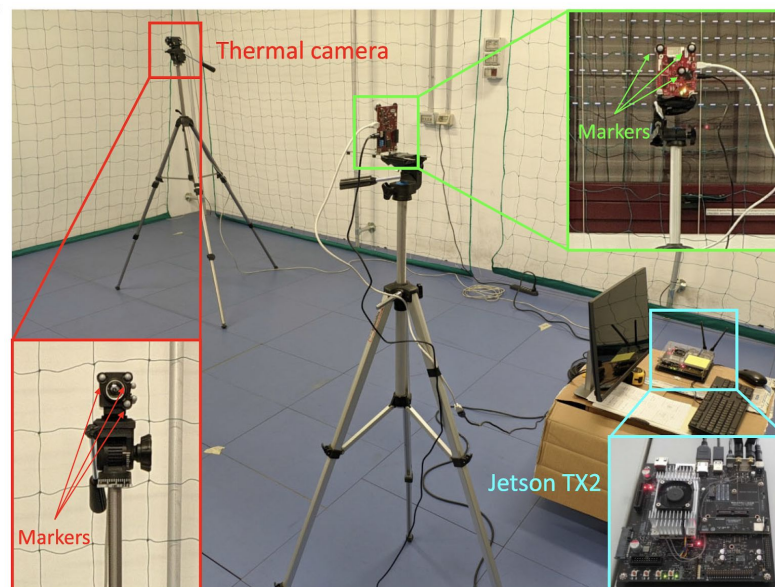
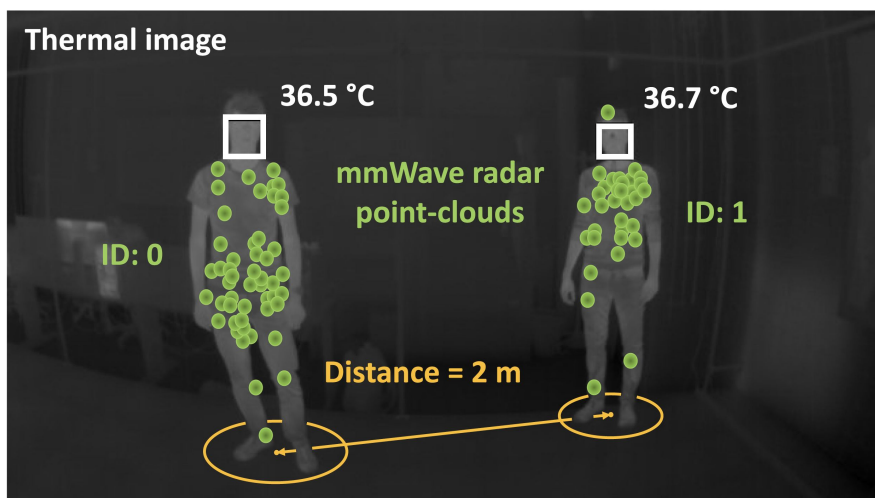
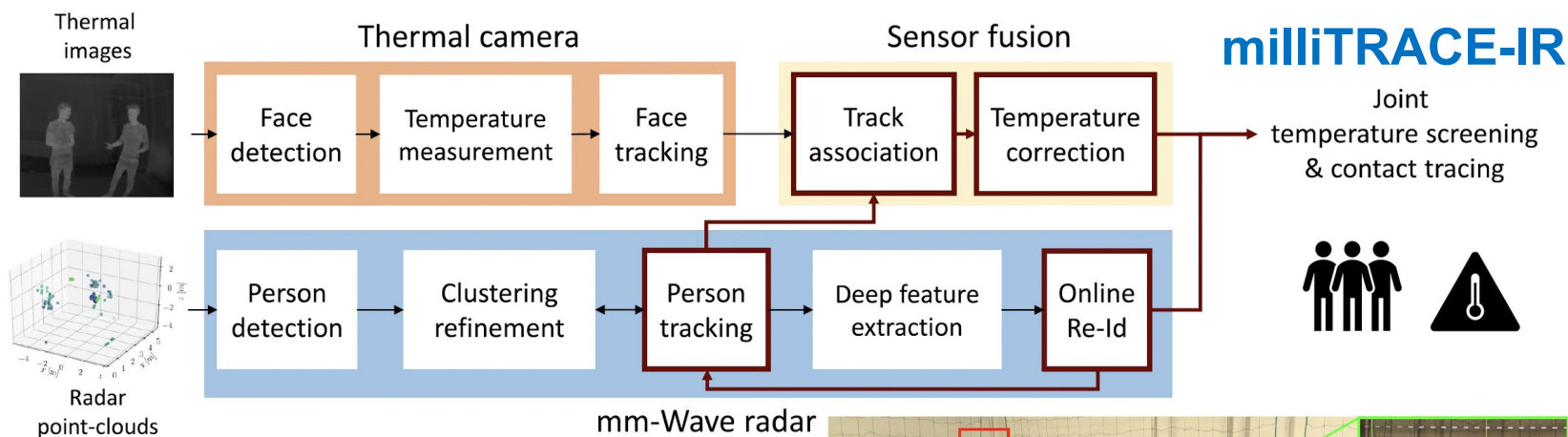


Identification



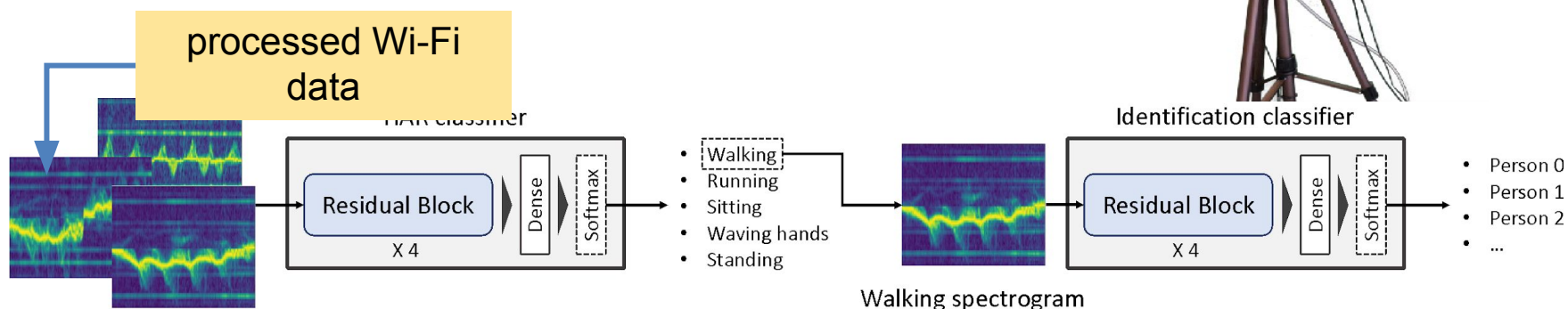
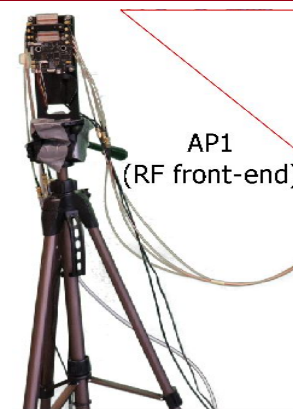
Radar position

Temperature and contact tracing

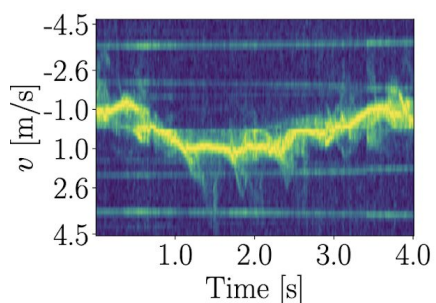


Wi-Fi sensing

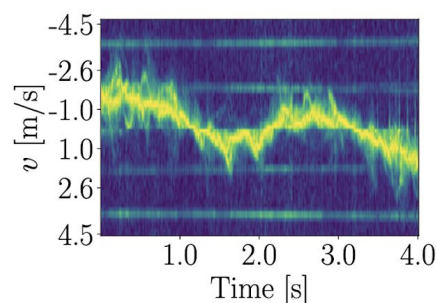
RAPID – indoor human detection and sensing through a testbed implementing the new IEEE 802.11ax Wi-Fi standard at 60 GHz



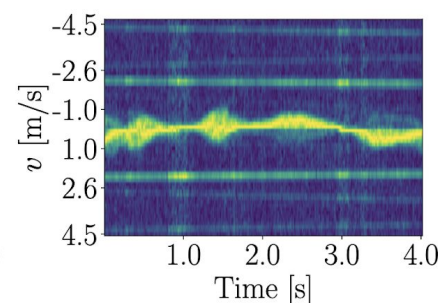
micro-Doppler input



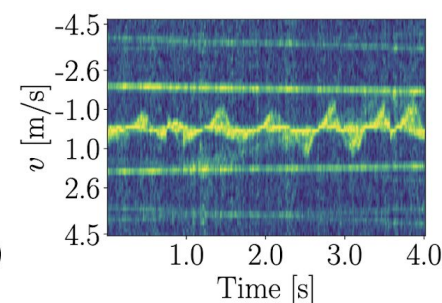
Walking.



Running.



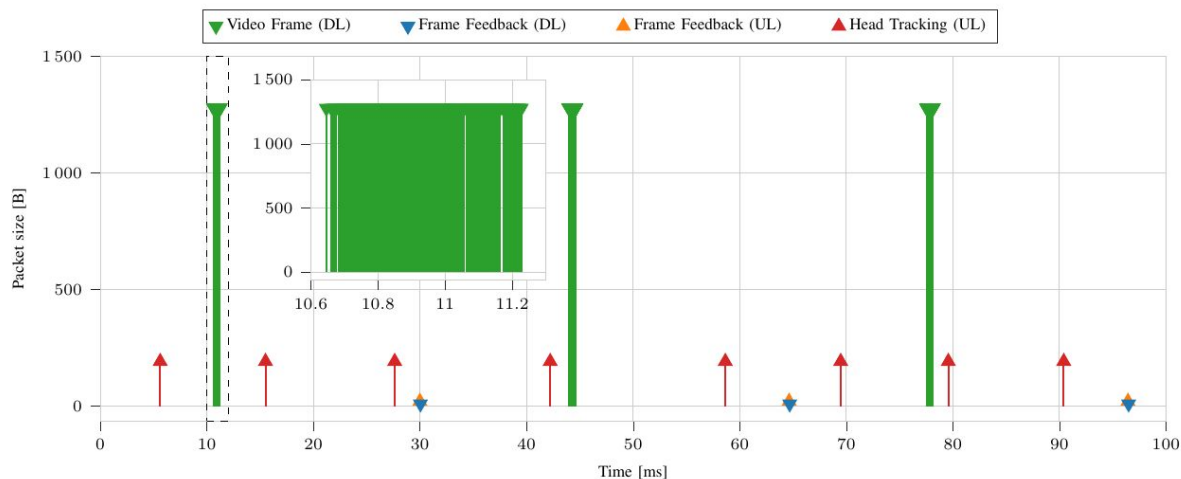
Sitting down.



Waving hands.

VR trace analysis

- ❑ How do we predict VR traffic?
- ❑ Frames depend on activity: what is the user doing?
- ❑ ML applied on traffic traces: capture and analysis tools



Semantic communications

- ❑ Adapting communications to only send the most relevant information
- ❑ Mix of ML styles: reinforcement, supervised, unsupervised
- ❑ Theory of mind: how do we model other agents?

