

Summarization in Sequential Data with Applications to Procedure Learning

Ehsan Elhamifar

**College of Computer and Information Science
Northeastern University**

Data deluge



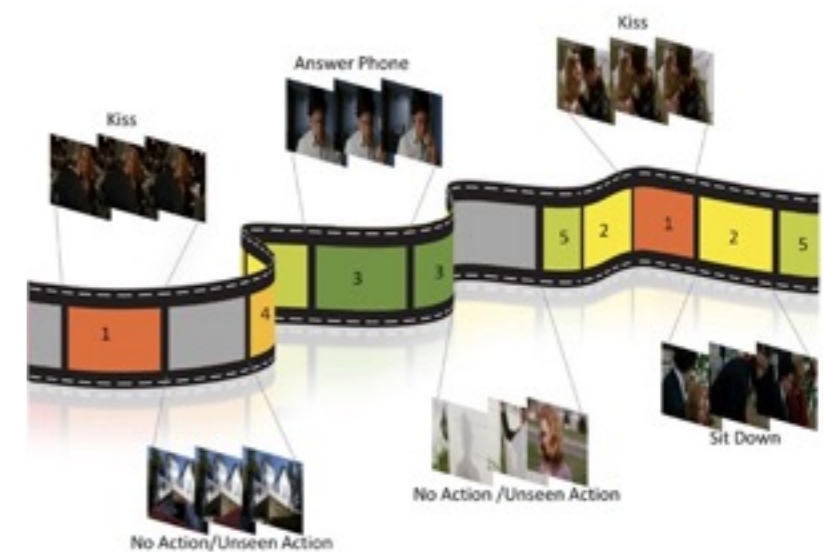
> 400hrs videos / minute



> 300M photos / day



24/7/365



Congratulations! Movies we think You will ❤️

Add movies to your Queue, or Rate ones you've seen for even better suggestions.

 Add ★★★★☆ Not Interested	 Add ★★★★☆ Not Interested	 Add ★★★★☆ Not Interested	 Add ★★★★☆ Not Interested
 Play Add ★★★★☆ Not Interested	 Add ★★★★☆ Not Interested	 Add ★★★★☆ Not Interested	 Play Add ★★★★☆ Not Interested

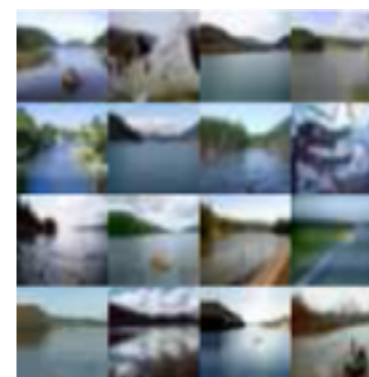
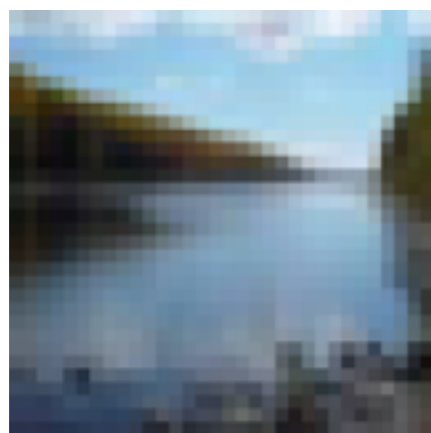
amazng awesome bad band beautiful berlin best better
book chocolate coffee come coming congrats cool creative
ch definitely dinner divx doing editing email end epic-fu
feeling feels finished food forward @frankjonen free frier
me getting girl giving glad goes going google gr
hah halo happy having head heading hear hearing ho
stsetshow.com ill internet iphone @irinaslutsky isnt jetse
left life listening little live lol long look looking looks
ke makes making man maybe media missed morning m
nice night nyc obama old omg online panel park party
phone pixelodeon play playing post pretty reading ri
ey right rock run said saw say says @seanbonner serio
is @spin @spytap start steve @stevewoolf stop strik
ing teh tell thanks @thefemgeek theyre thing things t
t thx time tiny today tomorrow tonight @tonykatz try
!O vote wait waiting want watch watching web week w
wondering work working world writers writing wrong yay

Big data opportunities

- Capturing modes and patterns **unseen** in small data
 - Increase learning and inference performance
- From **Nearest Neighbors**



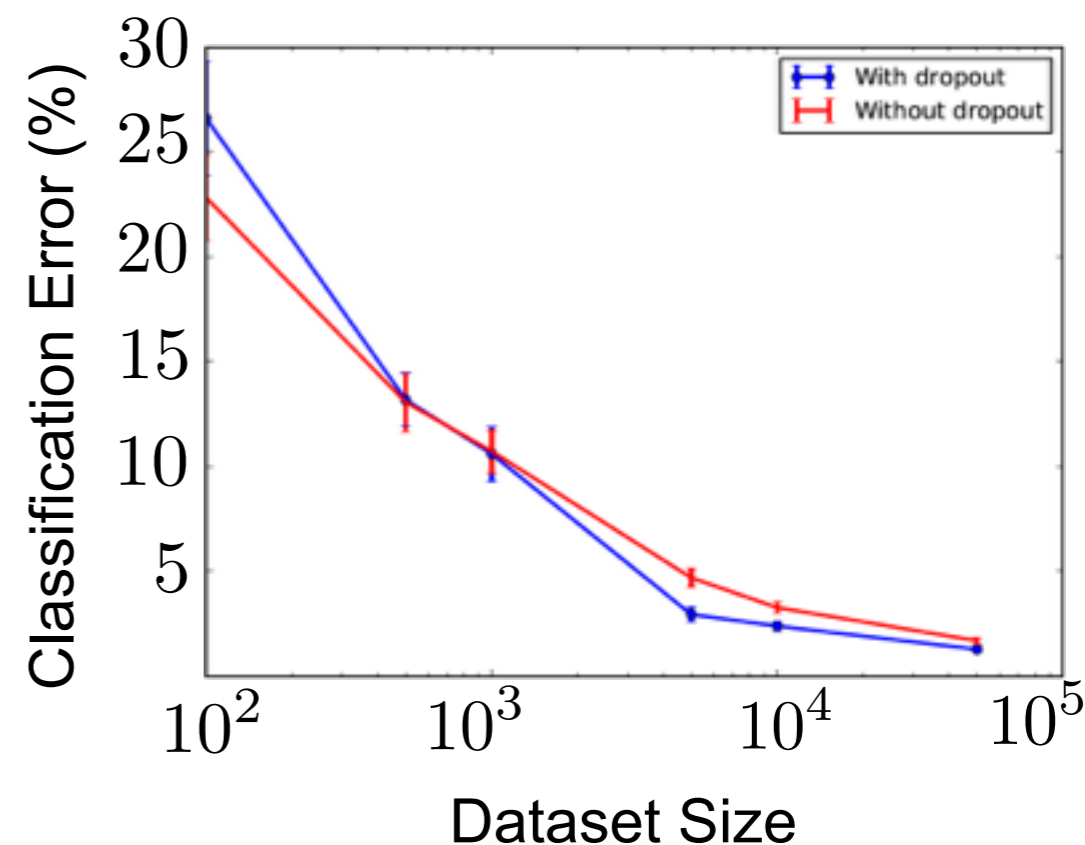
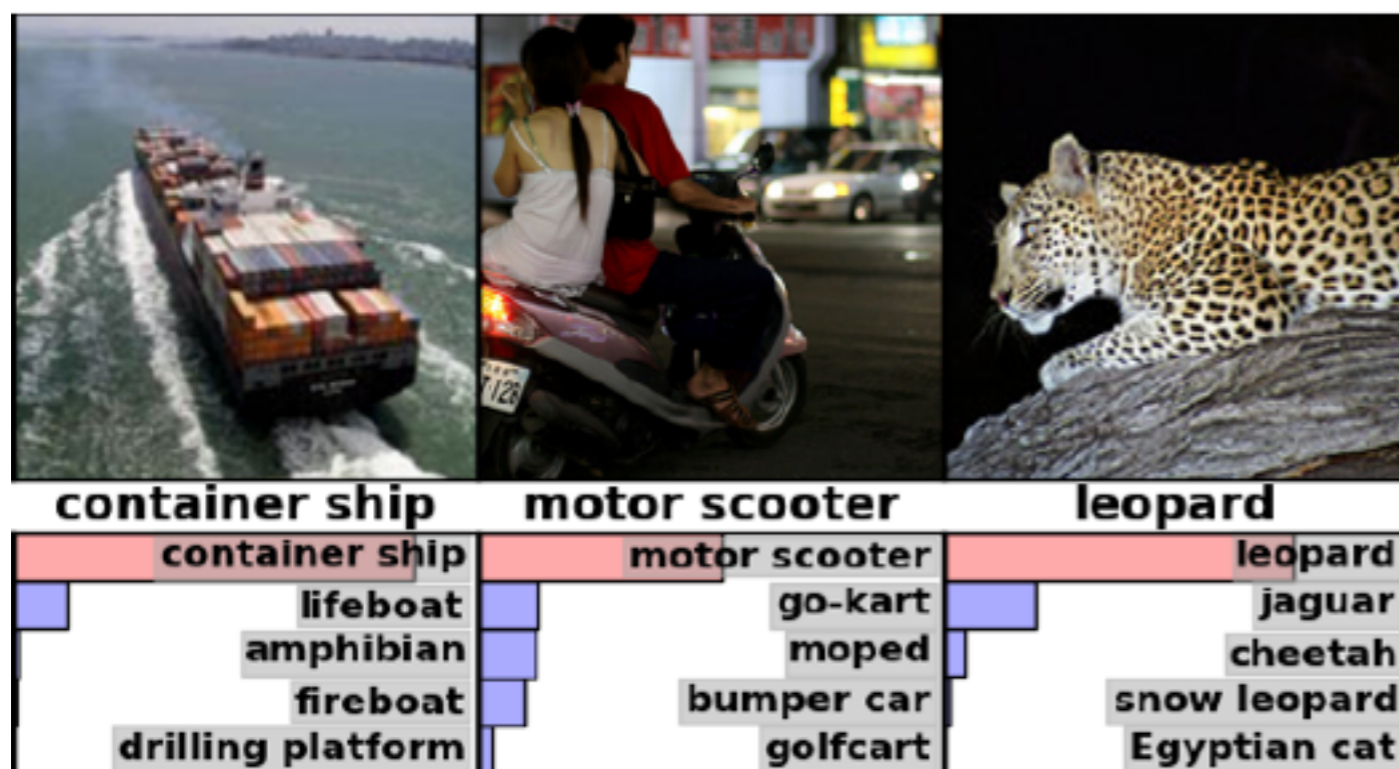
“It takes a large amount of data for our method to succeed. We saw **dramatic improvement** when moving from ten thousand to two million images.” [Hayes, Efros, SIGGRAPH '07]



“As we increase the size of the dataset from 10^5 to the 10^8 images, the quality of the retrieved set **increases dramatically.**” [Torralba, Fergus, Freeman, PAMI '08]

Big data opportunities

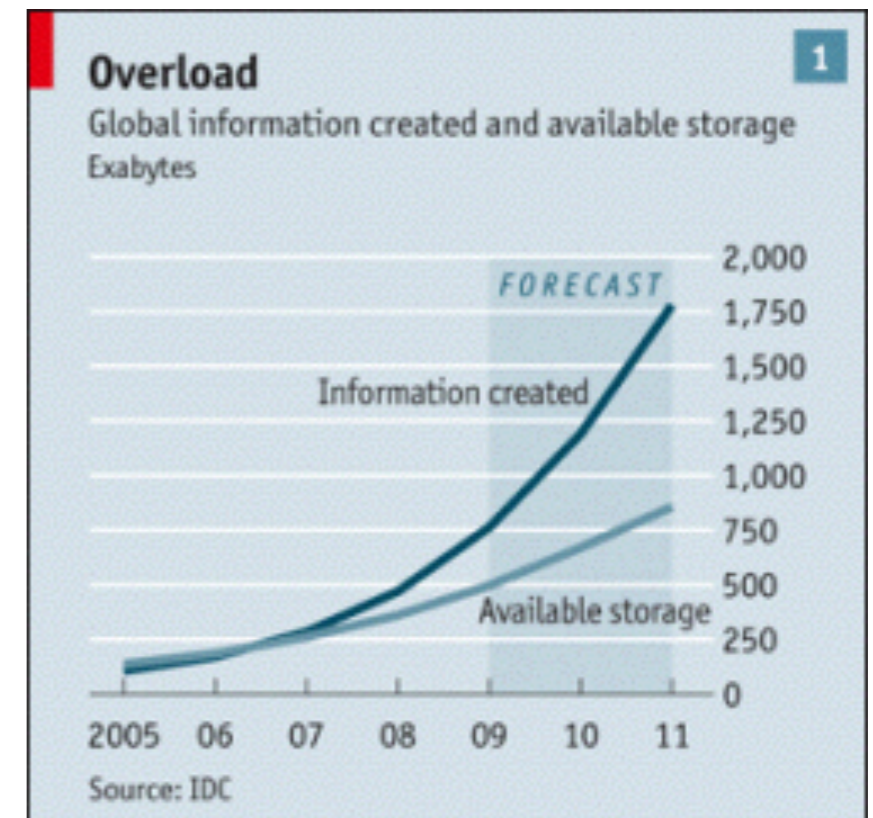
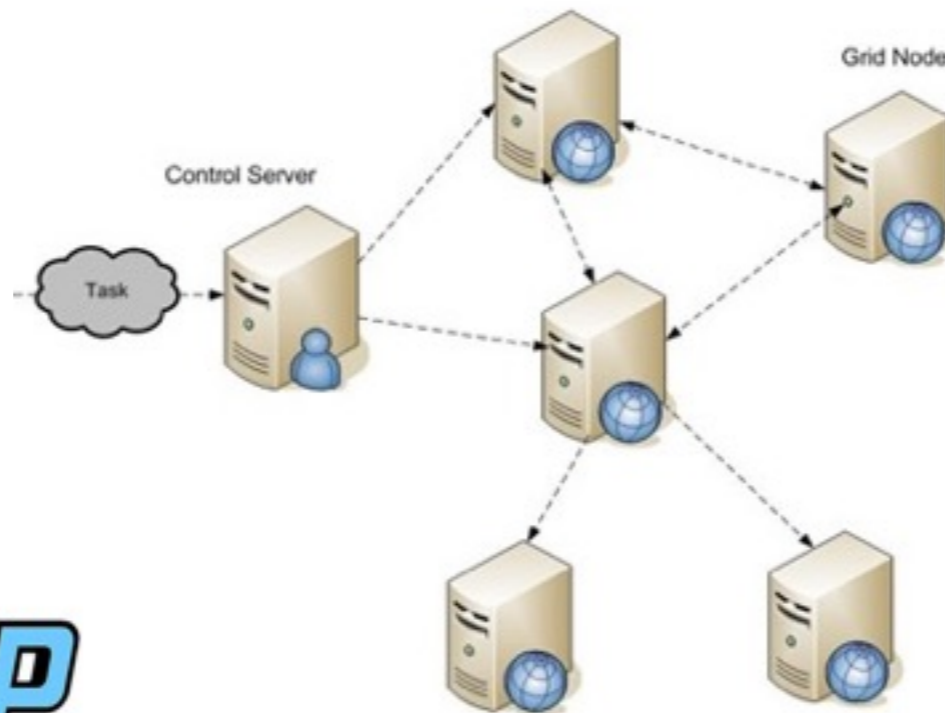
- Capturing modes and patterns **unseen** in small data
 - Increase learning and inference performance
- To **Deep Neural Networks**



“From 100 to 100,000 training samples, classification error **drops from 25% to less than 5%.**” [Srivastava-Hinton-Krizhevsky-Sutskever-Salakhutdinov '14]

Big data challenges

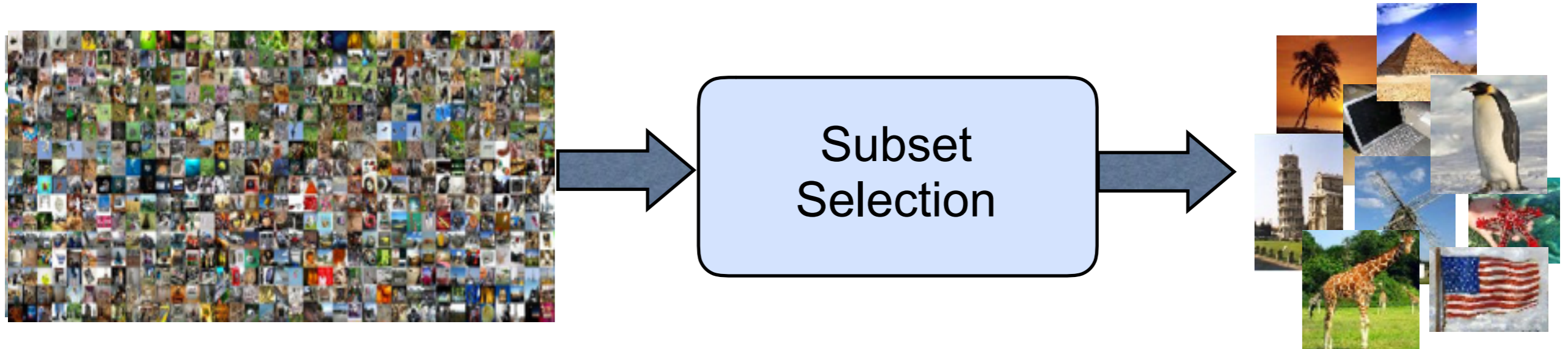
- Capturing modes and patterns **unseen** in small data
 - Increase learning and inference performance: from NNs to DNNs
- Limiting **computational** and **memory** resources
 - We have **more data** than capacities of our resources



Economist, Feb 2010

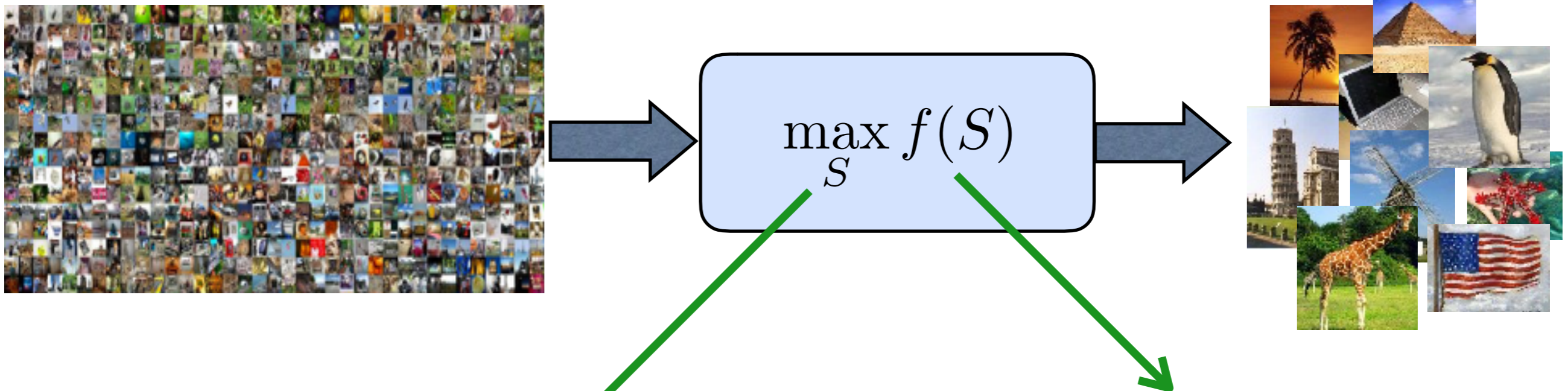
Subset selection

- Find **small data** capturing **statistical properties** of large data



Subset selection components

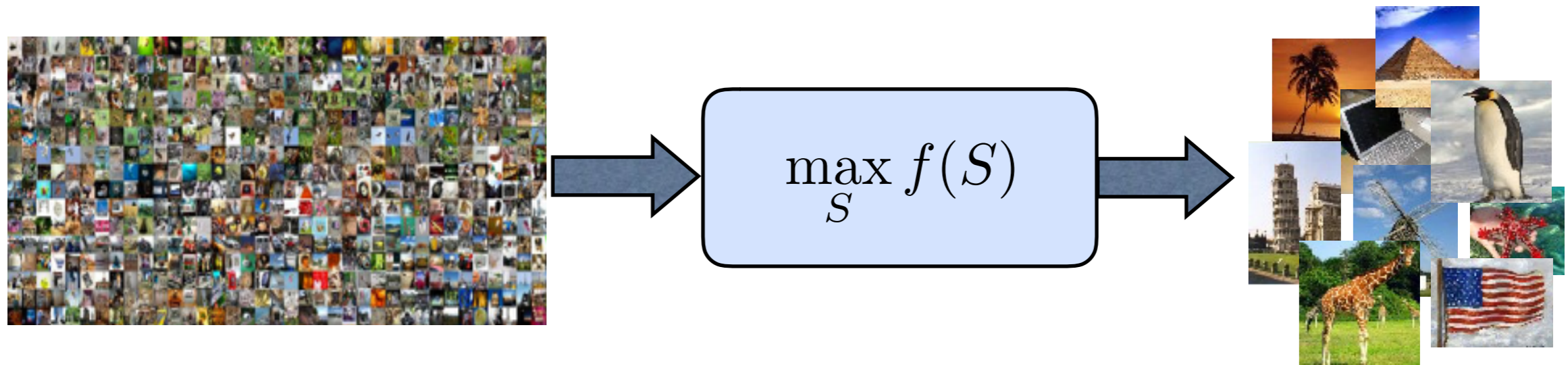
- How to **characterize informativeness** of items?
- How to recover the **most informative** items?



2. How to optimize?

1. What objective functions?

This Talk



- Subset selection is **NP-hard**, hence, **approximate** algorithms
 - **Better summarization**: closer to **global optimum**
- Focus on one class of objective functions: **Facility Location**
 - ➔ Develop **efficient sparse optimization**
 - ➔ Show it can **achieve global optimum** (under certain condition)

This Talk

- **Sequential data** (video, text, signals,...): large part of **‘Big Data’**

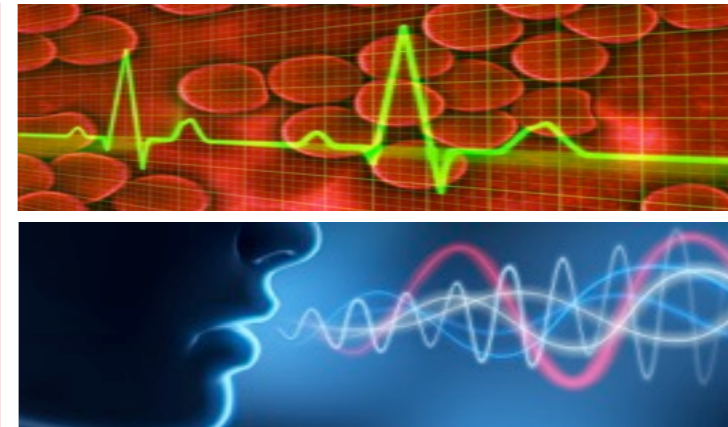


(CNN) -- A Japanese rocket roared into orbit early Friday (Thursday afternoon ET) carrying what NASA calls its most precise instrument yet for measuring rain and snowfall.

The Global Precipitation Measurement (GPM) satellite is the first of five earth science launches NASA has planned for 2014. The 4-ton spacecraft is the most sophisticated platform yet for measuring rainfall, capable of recording amounts as small as a hundredth of an inch an hour, said Gail Skofronick Jackson, GPM's deputy project scientist.

The \$900 million satellite is a joint project with the Japanese space agency JAXA, and it lifted off from Tanegashima Space Center at 3:37 a.m. Friday (1:37 p.m. Thursday ET). In a little over a half hour, it had reached orbit, deployed its solar panels and began beaming signals back to its controllers, NASA said.

Also, once fully activated, GPM will use both radar and microwave instruments to detect falling snow for the first time. It will also combine data from other satellites with its own readings, beaming



- Sequential data have **structured dependencies**
 - **Instructional videos**: ordering of key steps

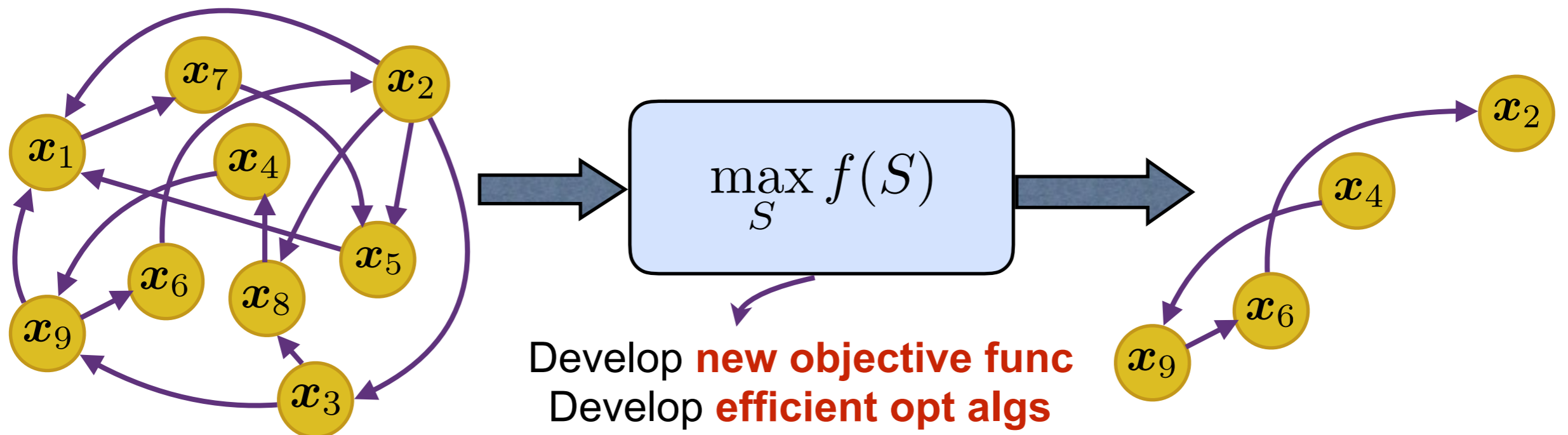


Key steps



This Talk

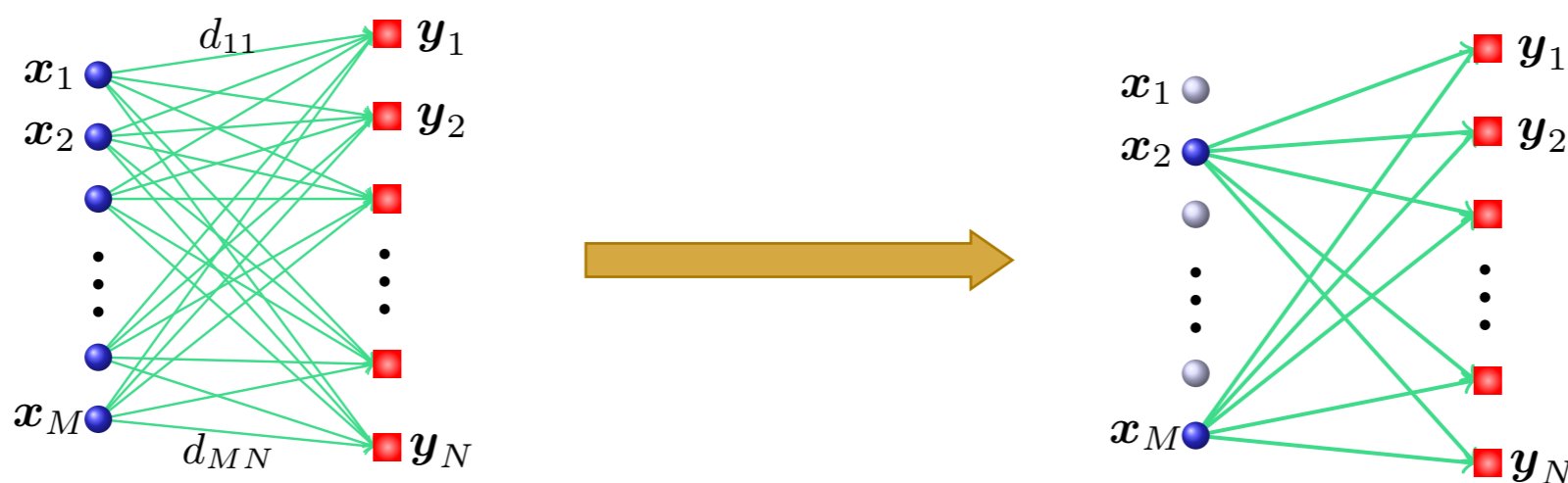
- **Sequential data** (video, text, signals,...): large part of **'Big Data'**
- Sequential data have **structured dependencies**
 - Existing SS: data is a **bag of randomly permutable** items
- **Generalize facility location** to handle **seq summarization**



Subset selection from dissimilarities

- Given dissimilarities $d : \mathbb{X} \times \mathbb{Y} \longrightarrow \mathbb{R}^{\geq 0}$,
source target

select **a small subset** of \mathbb{X} that **well represent** \mathbb{Y} w.r.t. $d(\cdot, \cdot)$

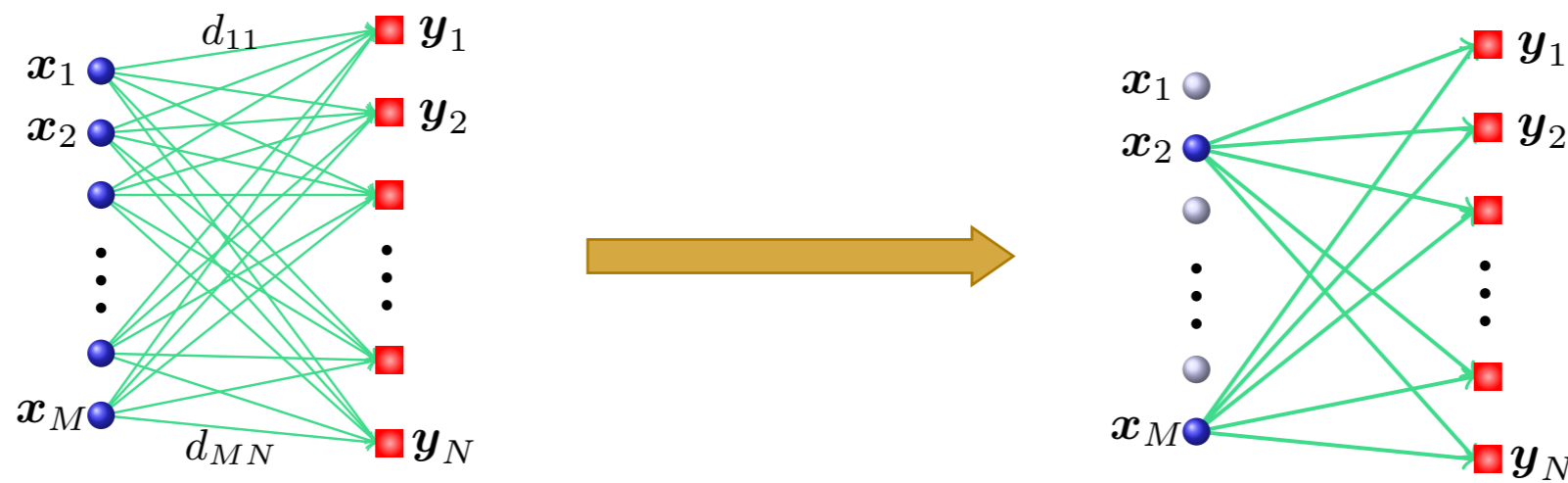


- $d(\mathbf{x}_i, \mathbf{y}_j) = d_{ij}$: how well \mathbf{x}_i represents \mathbf{y}_j
 - $\mathbb{X} = \text{data}, \quad \mathbb{Y} = \text{data} \longrightarrow d(\cdot, \cdot) = \text{Euclidean/geodesic distance}$
 - $\mathbb{X} = \text{models}, \quad \mathbb{Y} = \text{data} \longrightarrow d(\cdot, \cdot) = \text{representation/coding error}$

Subset selection from dissimilarities

- Given dissimilarities $d : \mathbb{X} \times \mathbb{Y} \longrightarrow \mathbb{R}^{\geq 0}$,
 \downarrow \downarrow
source target

select **a small subset** of \mathbb{X} that **well represent** \mathbb{Y} w.r.t. $d(\cdot, \cdot)$



- $d(\mathbf{x}_i, \mathbf{y}_j) = d_{ij}$: how well \mathbf{x}_i represents \mathbf{y}_j

Not necessarily metric:
asymmetric,
violate triangle ineq



Dissimilarity-based Sparse Subset Selection

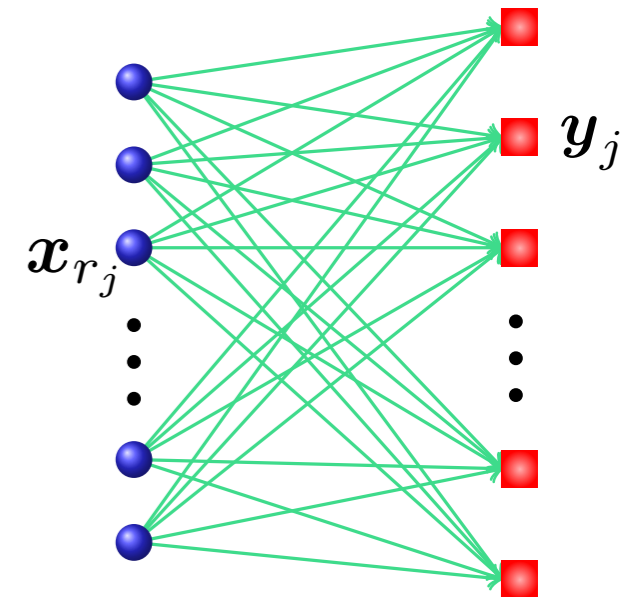
- Let $r_j \in \{1, \dots, M\}$ be index of representative of y_j
- **Approach:** define & maximize **potential function**

$$\Psi(r_1, \dots, r_N) \triangleq \Phi_{\text{enc}}(r_1, \dots, r_N) \times \Phi_{\text{card}}(r_1, \dots, r_N)$$

$$\Phi_{\text{enc}}(r_1, \dots, r_N) \triangleq \prod_{j=1}^N e^{-d_{r_j, j}} \quad \longrightarrow \quad \text{Prefer lower encoding}$$

$$\Phi_{\text{card}}(r_1, \dots, r_N) \triangleq e^{-\lambda |\{r_1, \dots, r_N\}|} \quad \longrightarrow \quad \text{Prefer small \# representatives}$$

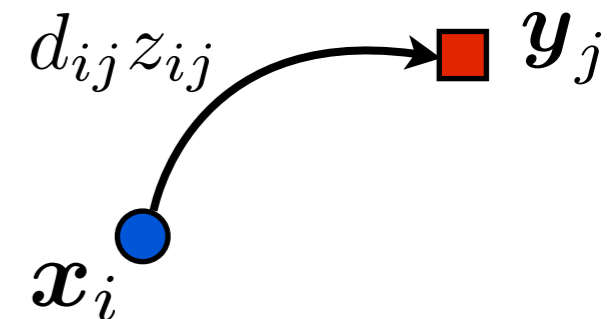
- Max Ψ over $\{r_1, \dots, r_N\} \in \{1, \dots, M\}^N$: **combinatorial**



DS3 formulation

- Develop a **sparse binary optimization**

- Define variables $z_{i,j} \triangleq \mathbb{I}(r_j = i) \in \{0, 1\}$



$$\Psi = \prod_{i=1}^M \prod_{j=1}^N \exp(-d_{i,j} z_{i,j}) \times \exp\left(-\lambda \sum_{i=1}^M \mathbb{I}(\| [z_{i,1} \dots z_{i,N}] \|_p)\right)$$

- Minimizing $-\log \Psi$:

$$\min_{\{z_{i,j}\}} \sum_{i=1}^M \sum_{j=1}^N d_{i,j} z_{i,j} + \lambda \sum_{i=1}^M \mathbb{I}(\| [z_{i,1} \dots z_{i,N}] \|_p)$$

Nonconvex

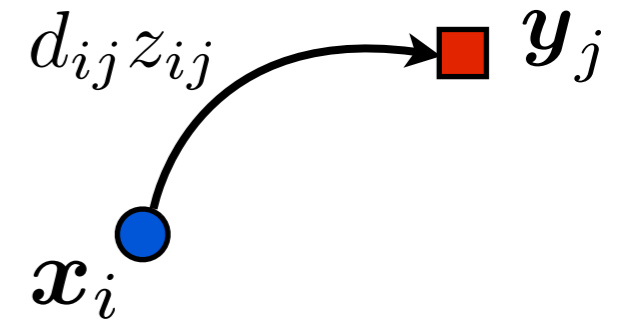
$$\text{s. t. } \underline{z_{i,j} \in \{0, 1\}}, \quad \sum_{i=1}^M z_{i,j} = 1, \quad \forall i, j$$

~ Facility location

DS3 formulation

- Develop a **sparse binary optimization**

- Define variables $z_{i,j} \triangleq \mathbb{I}(r_j = i) \in \{0, 1\}$



$$\Psi = \prod_{i=1}^M \prod_{j=1}^N \exp(-d_{i,j} z_{i,j}) \times \exp\left(-\lambda \sum_{i=1}^M \mathbb{I}(\| [z_{i,1} \dots z_{i,N}] \|_p)\right)$$

- Minimizing $-\log \Psi$:

Convex

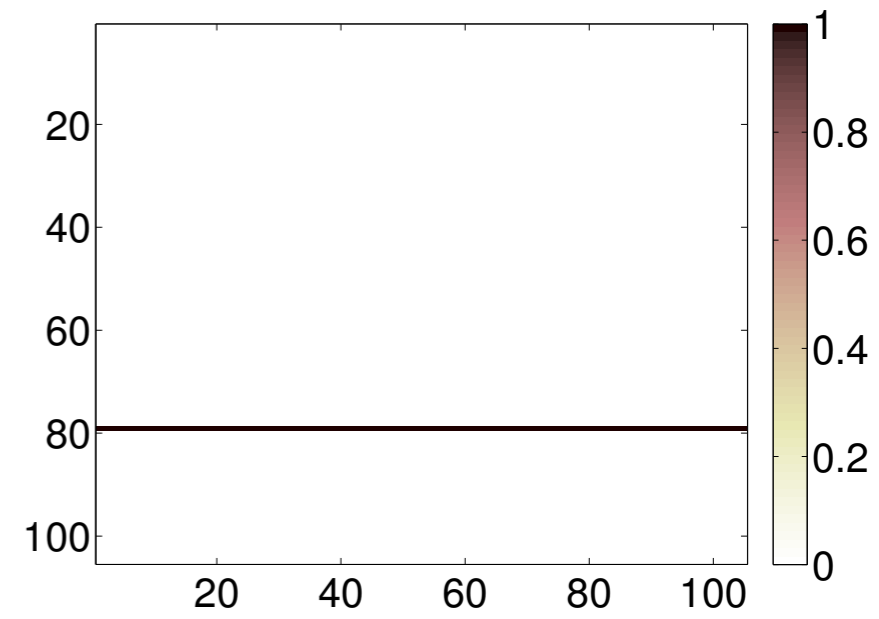
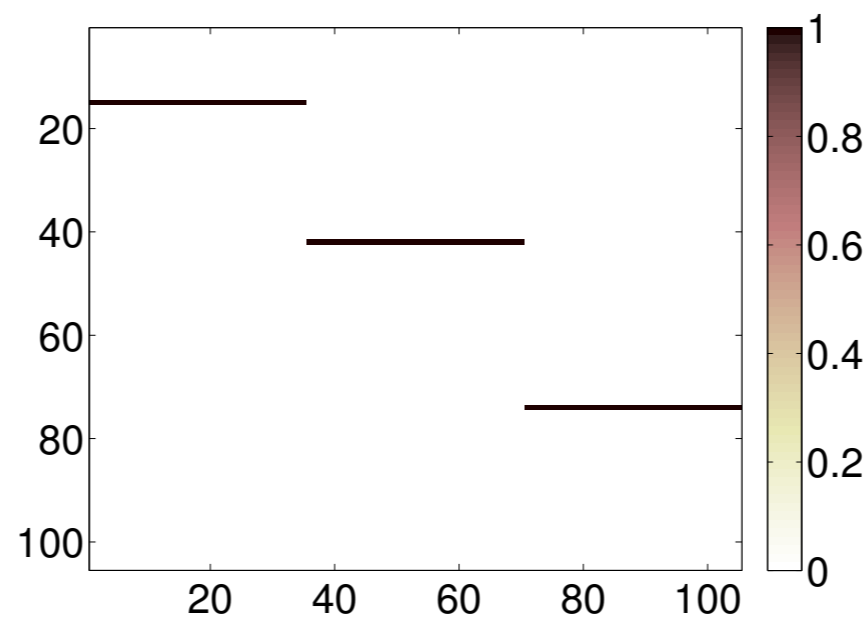
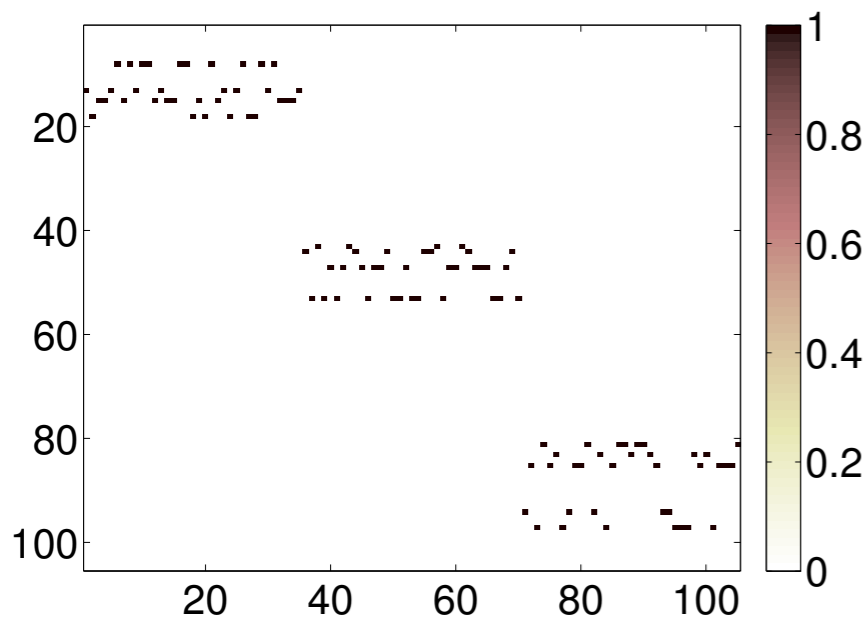
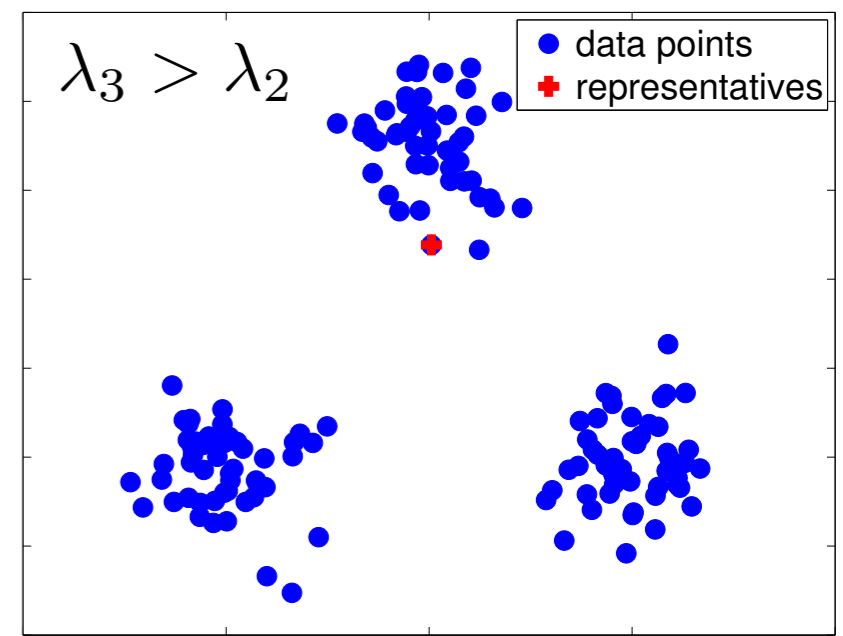
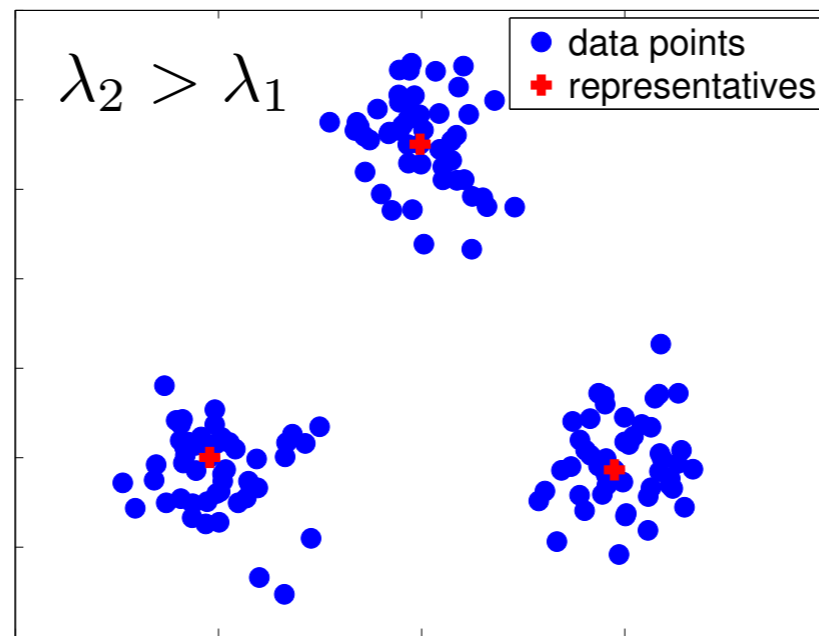
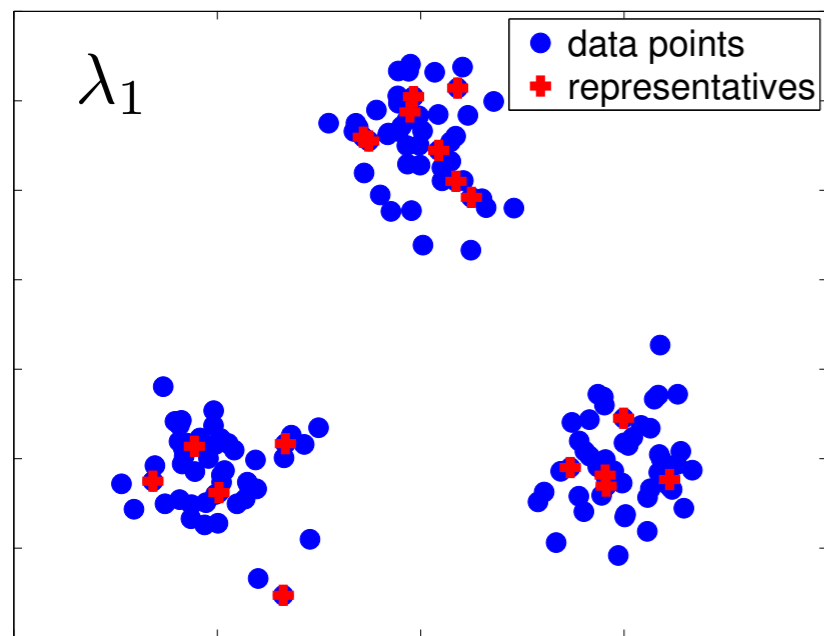
$$\min_{\{z_{i,j}\}} \sum_{i=1}^M \sum_{j=1}^N d_{i,j} z_{i,j} + \lambda \sum_{i=1}^M \underbrace{\| [z_{i,1} \dots z_{i,N}] \|_p}_{p \in \{2, \infty\}}$$

$$\text{s. t. } \underbrace{z_{i,j} \geq 0, \sum_{i=1}^M z_{i,j} = 1, \forall i, j}$$

DS3

Toy example

- Source = target = {data points}



$$\min_{\{z_{i,j}\}} \sum_{i=1}^M \sum_{j=1}^N d_{i,j} z_{i,j} + \lambda \sum_{i=1}^M \left\| [z_{i,1} \cdots z_{i,N}] \right\|_p \quad \text{s. t.} \quad z_{i,j} \geq 0, \quad \sum_{i=1}^M z_{i,j} = 1, \quad \forall i, j$$

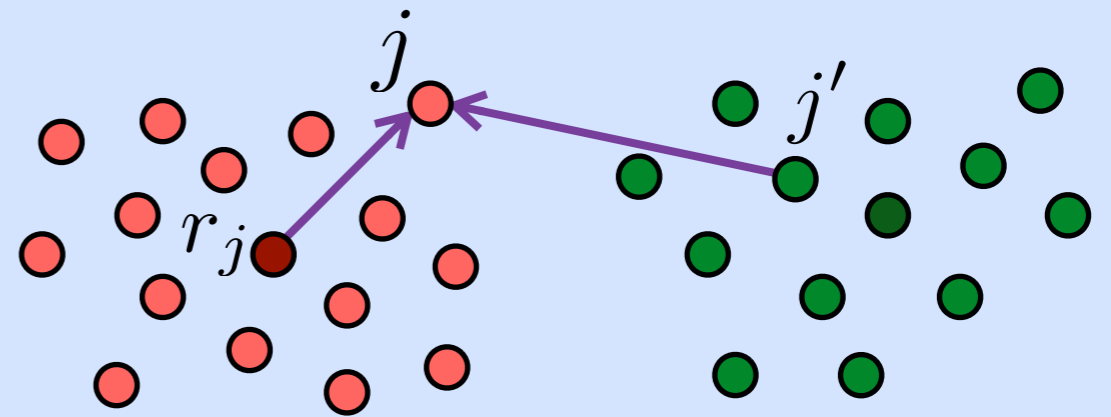
Theoretical analysis

- When is **convex relaxation exact**?

- **Theorem:** Let r_j be representative of j via **non-convex optimization**. The convex relaxation is exact if

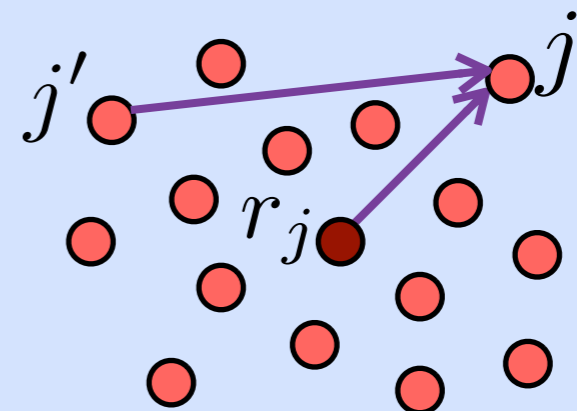
$$1) \quad \frac{\lambda}{N_{r_j}} + d_{r_j, j} < d_{j', j}$$

Clusters sufficiently separated



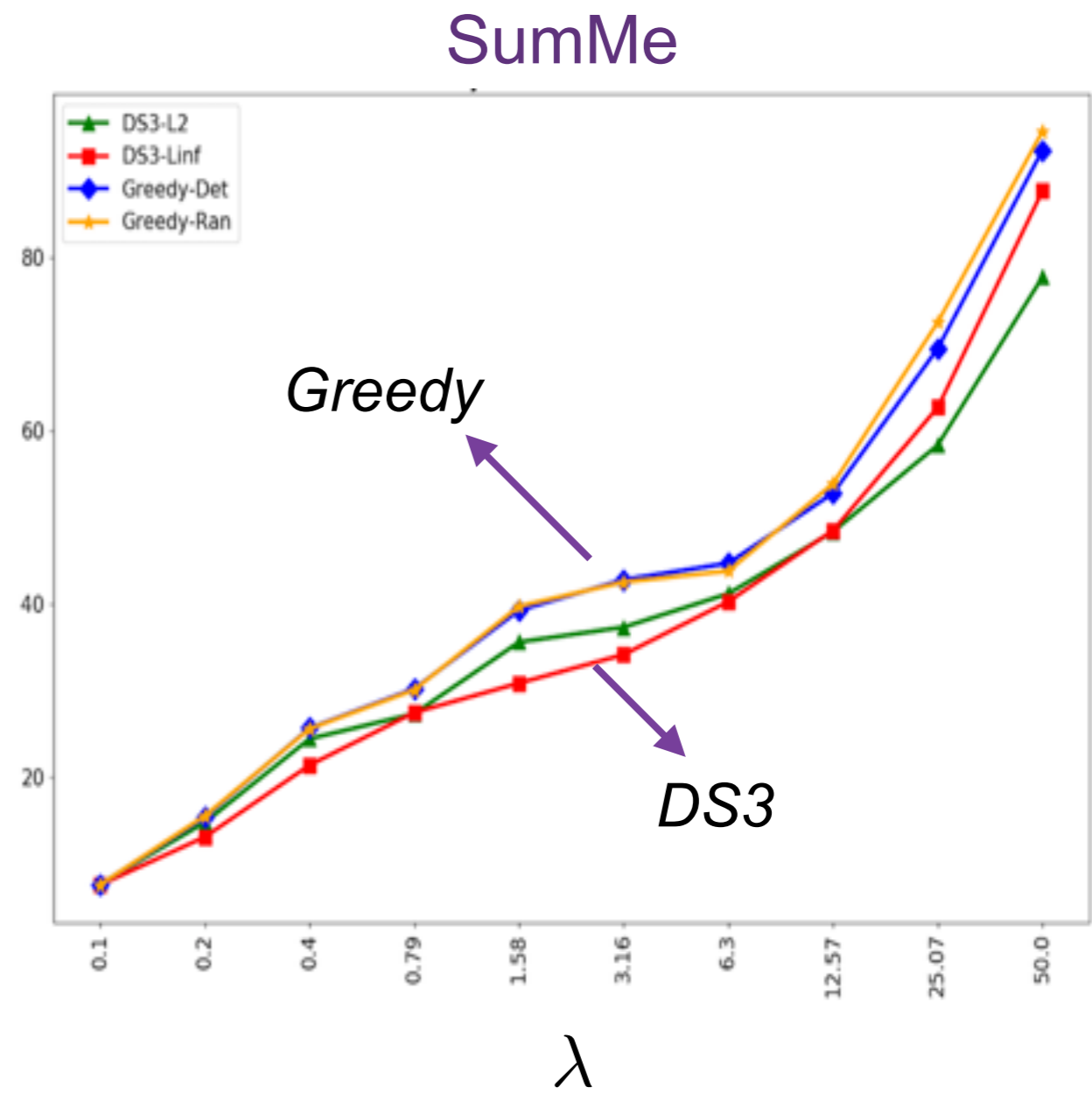
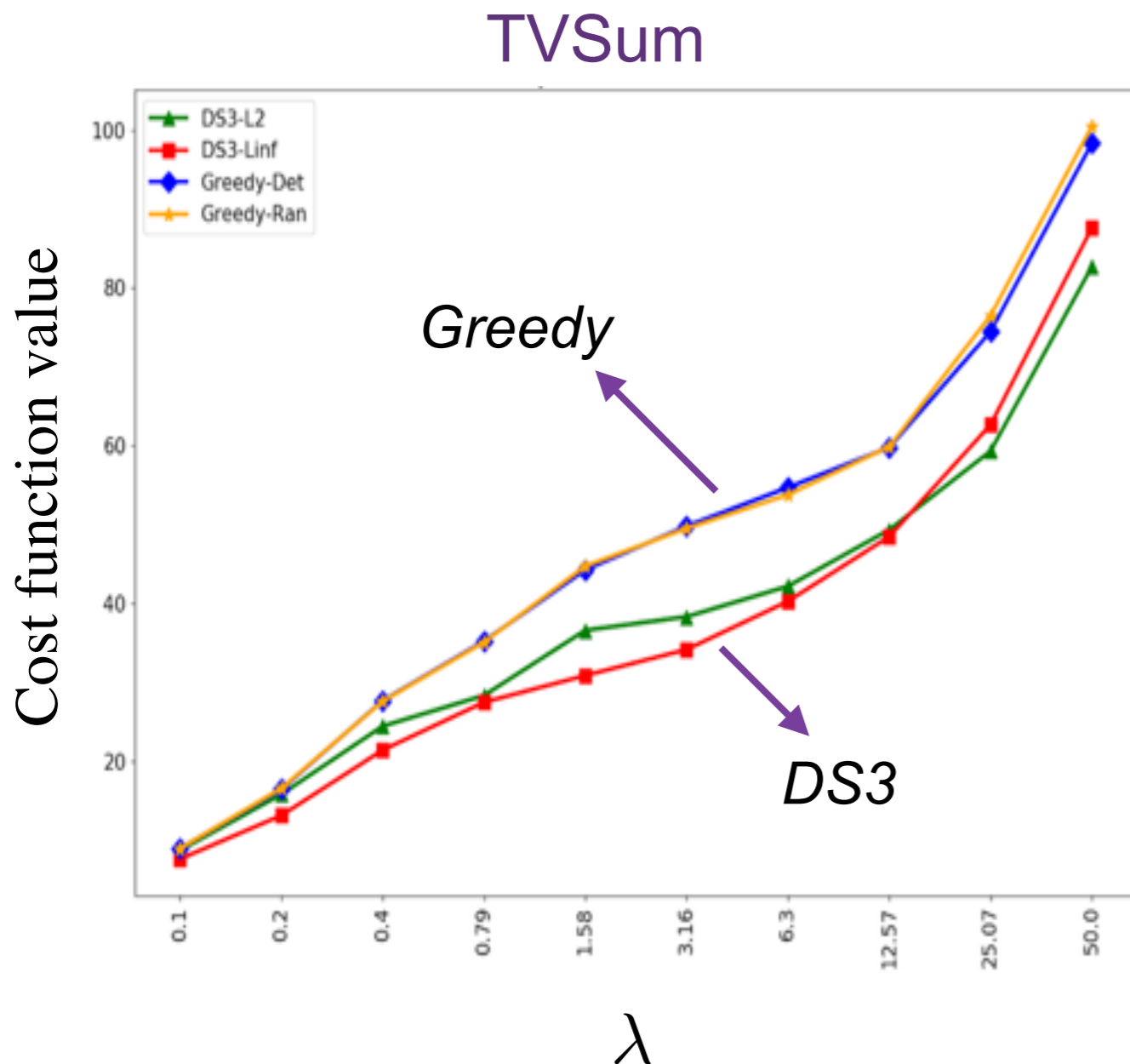
$$2) \quad d_{j', j} \leq \frac{\lambda}{N_{r_j}} + d_{r_j, j}$$

Sufficiently small radius



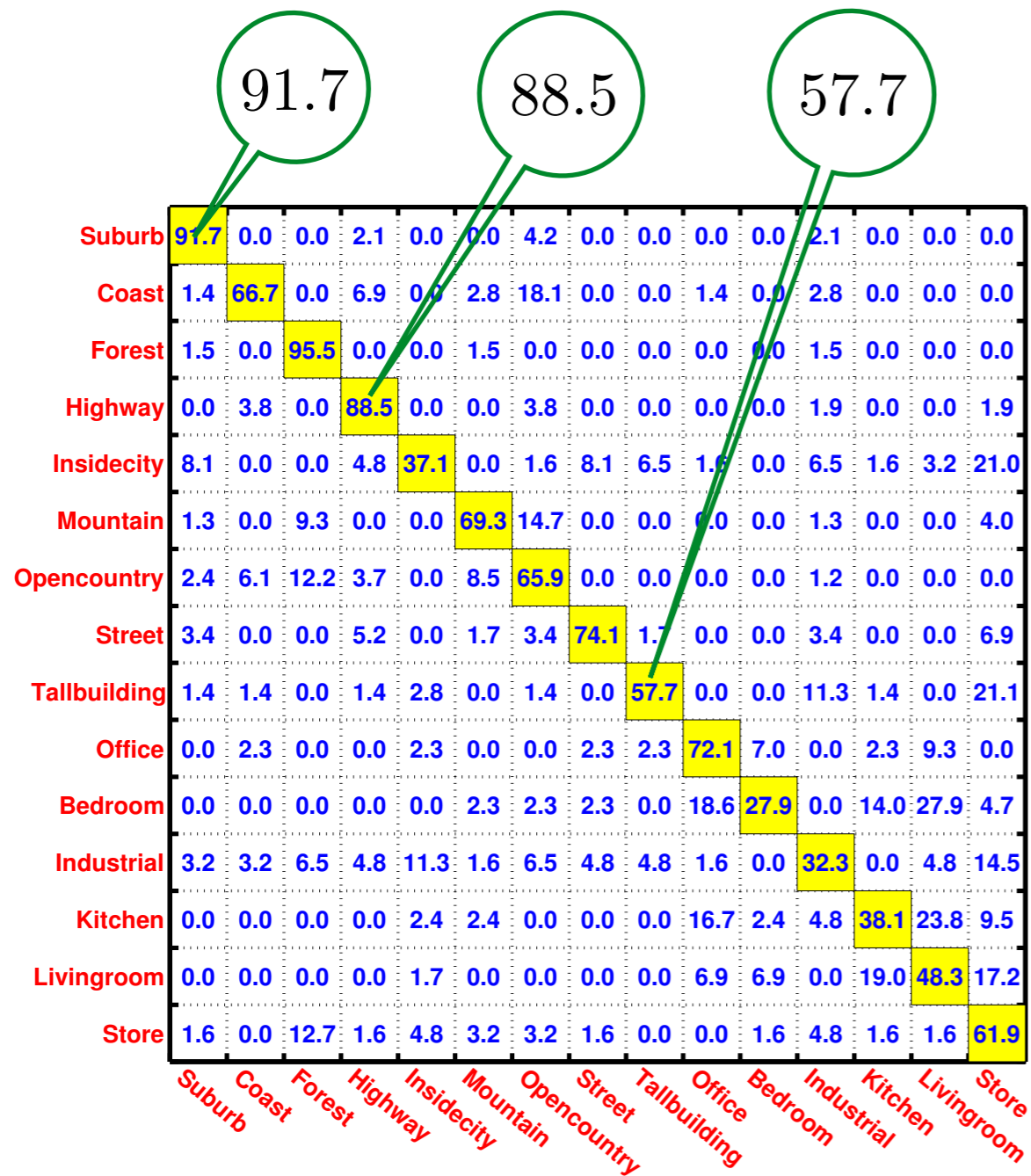
Comparison with greedy

- Video summarization
 - Solve FL via Greedy and sparse coding (DS3)

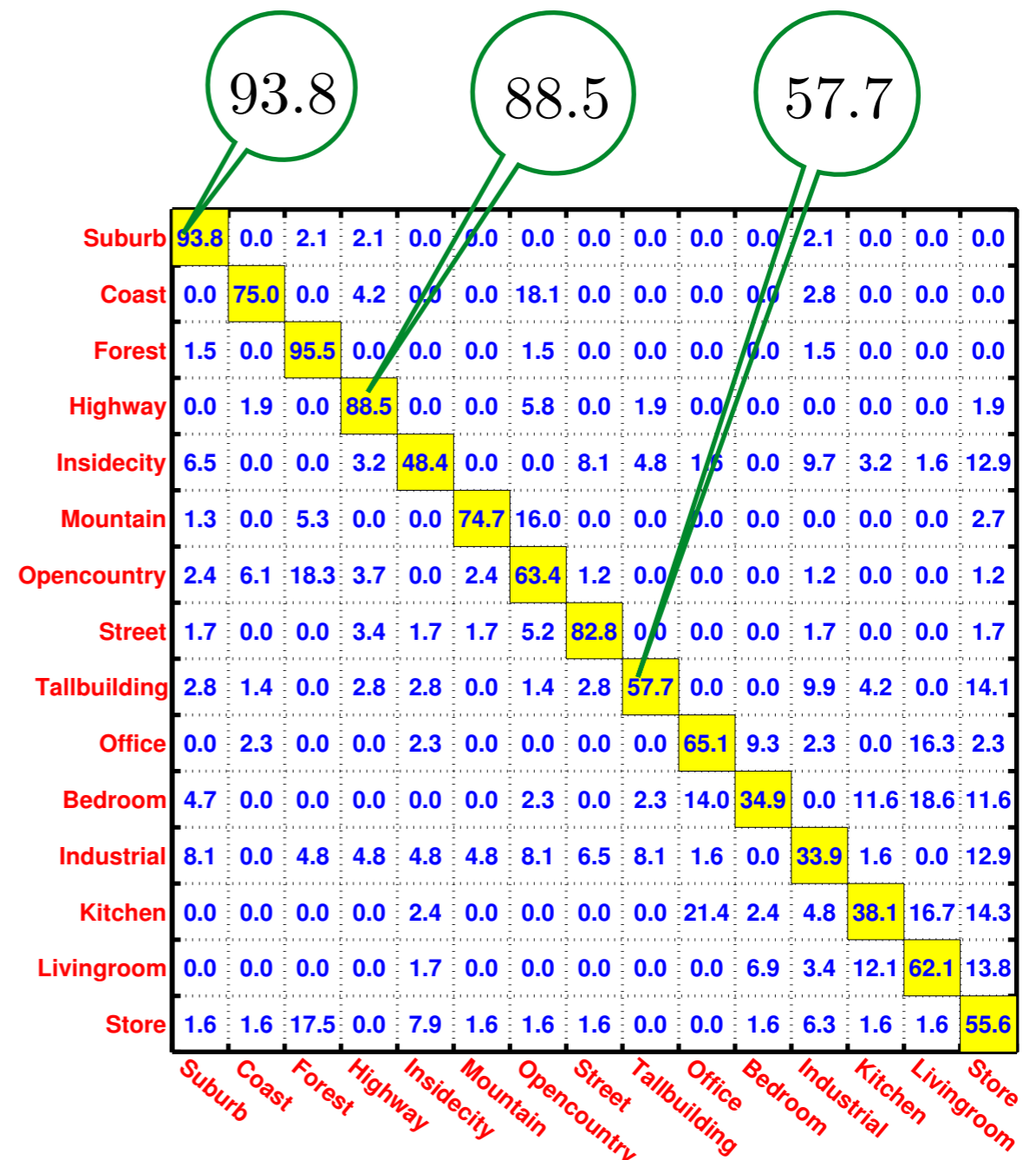


Recognition via representatives

- Select η fraction of training samples via DS3



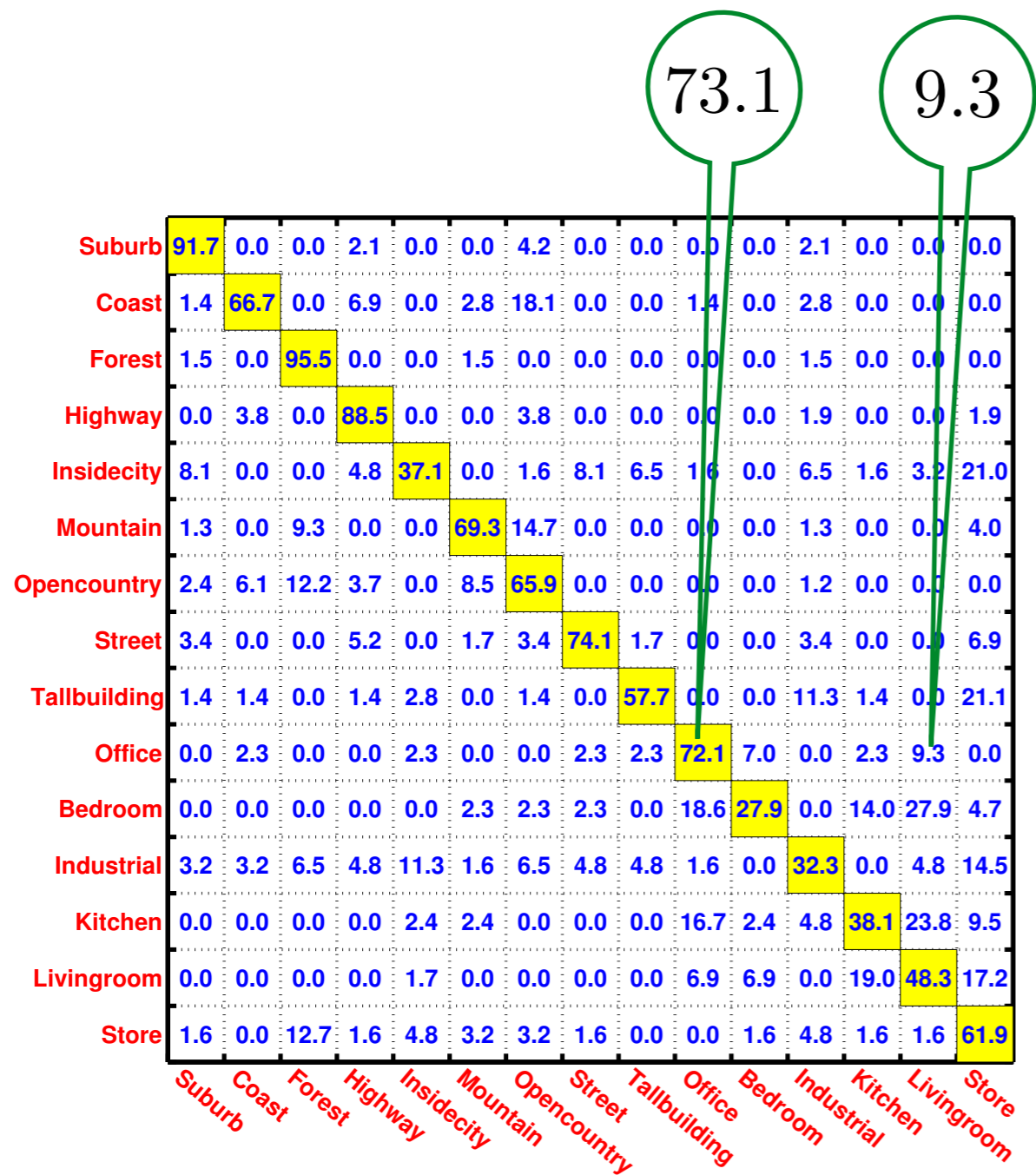
DS3 selecting 35% samples



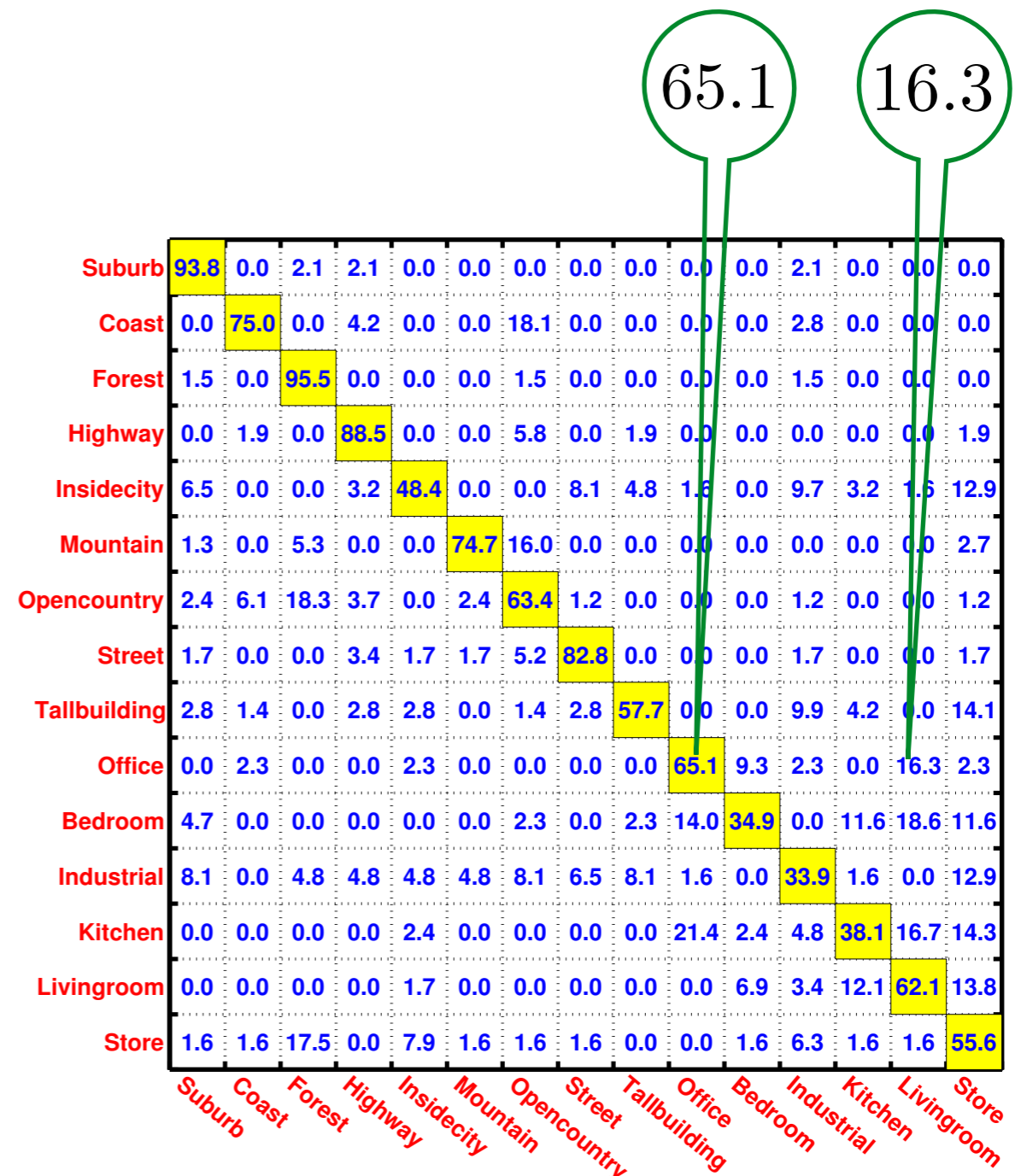
Trained via all samples

Recognition via representatives

- Select η fraction of training samples via DS3



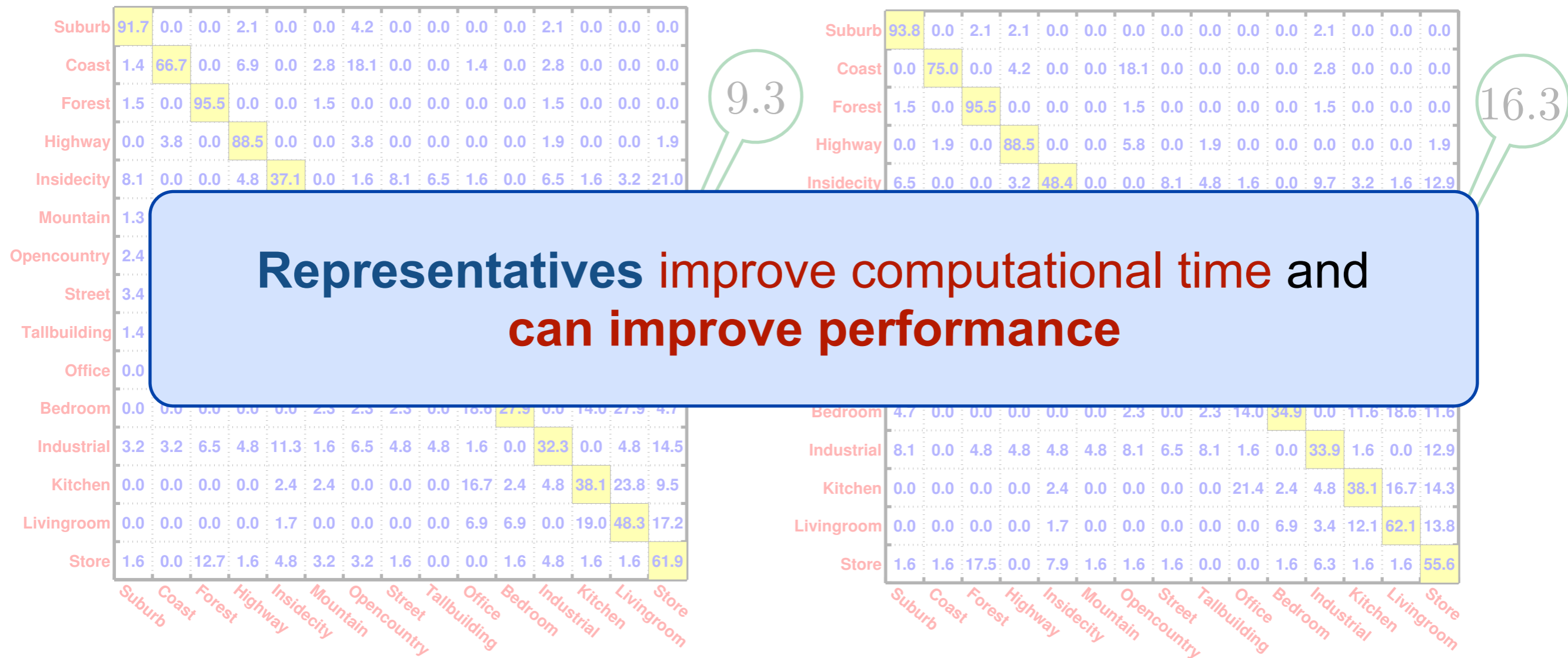
DS3 selecting 35% samples



Trained via all samples

Recognition via representatives

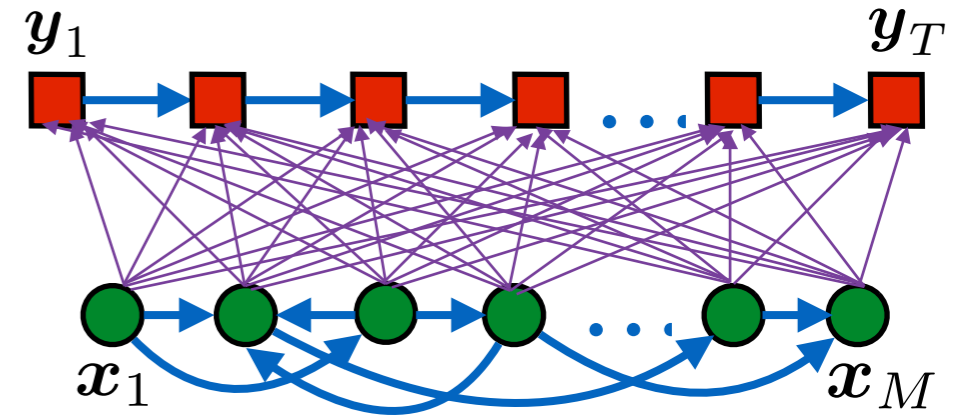
- Select η fraction of training samples via DS3



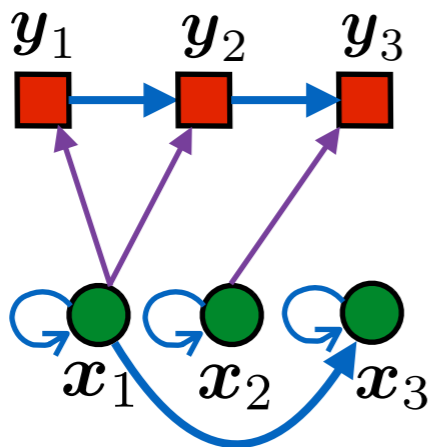
Sequential subset selection

- Incorporate **dynamics** of sequential data into SS

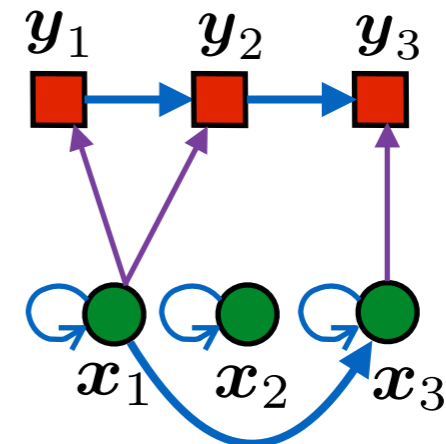
- Source set: with **transition dynamics**
- Target set: **sequential data**



- **Goal:** find (r_1, \dots, r_T) with **high transition probability** that **well encodes the data** and has **small cardinality**



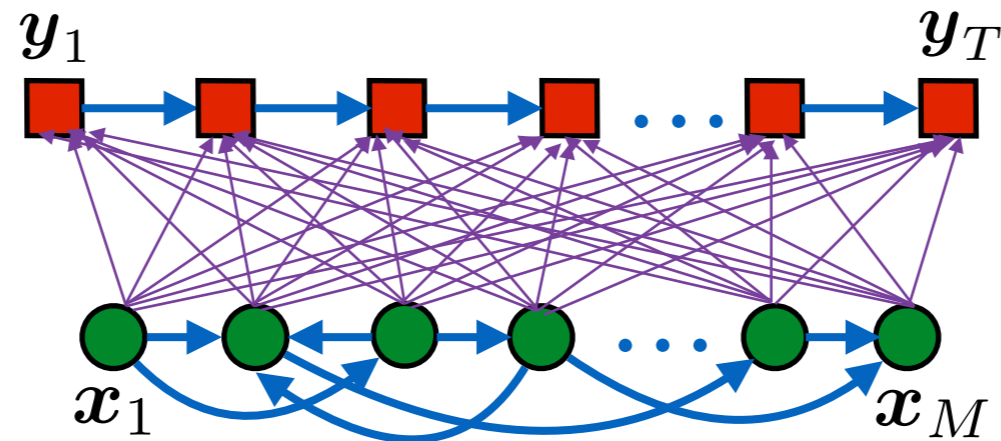
$$(r_1, r_2, r_3) = (1, 1, 2)$$



$$(r_1, r_2, r_3) = (1, 1, 3)$$

Sequential facility location

- Let r_t be index of representative of y_t



- Maximize global potential function

$$\max_{(r_1, \dots, r_T)} \Phi_{\text{enc}}(r_1, \dots, r_N) \times \Phi_{\text{card}}(r_1, \dots, r_N) \times \Phi_{\text{dyn}}(r_1, \dots, r_N)$$

Prefer a sequence of compatible reps

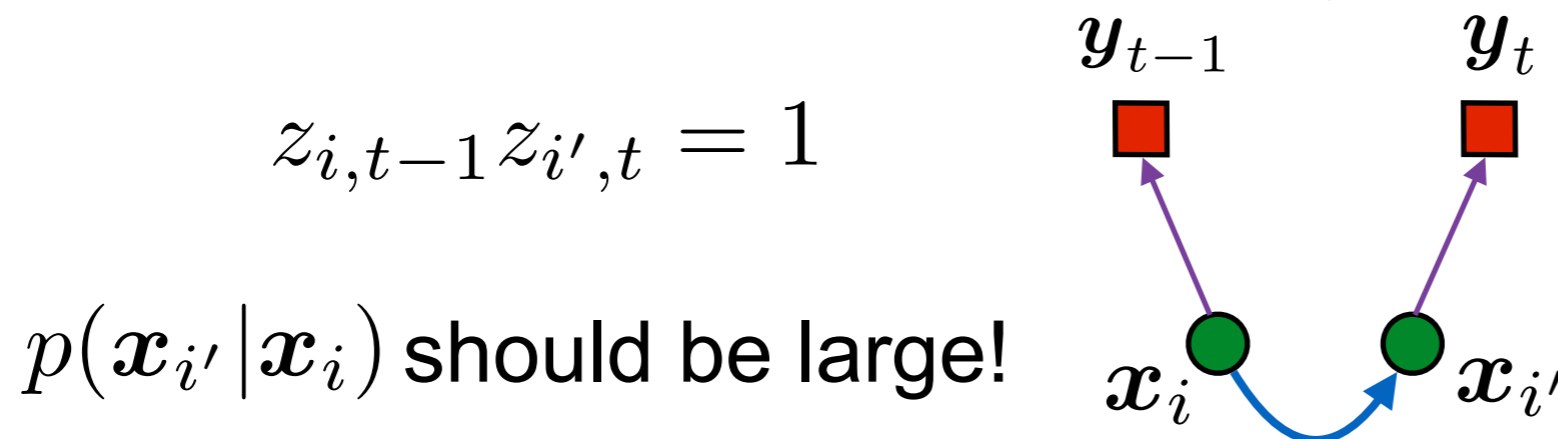
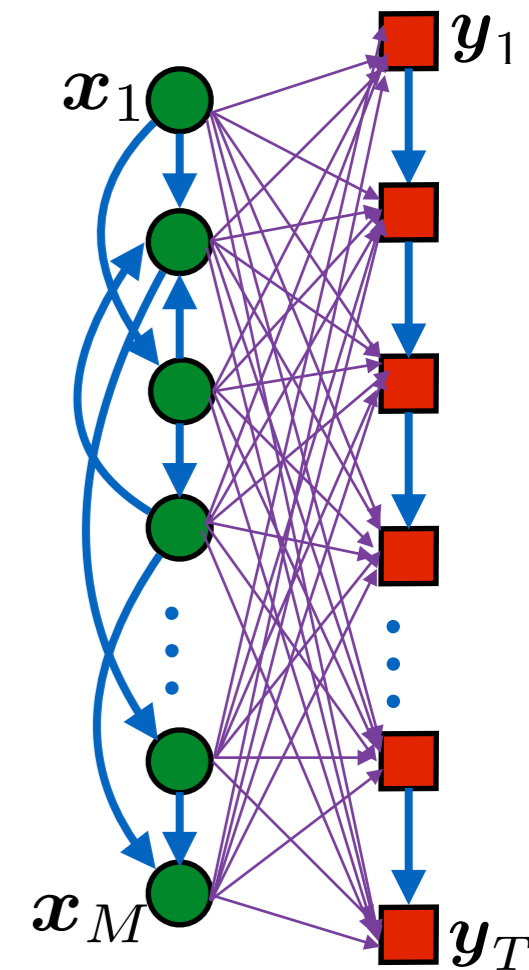
- **Dynamic potential:** n-th order Markov model

$$\Phi_{\text{dyn}}(r_1, \dots, r_N) \triangleq \left(\prod_t p_t(\mathbf{x}_{r_t} | \mathbf{x}_{r_{t-1}}, \dots, \mathbf{x}_{r_{t-n}}) \right)^\beta$$

Sequential facility location

- **SeqFL**: Solve optimization on assignment variables $\{z_{i,t}\}$

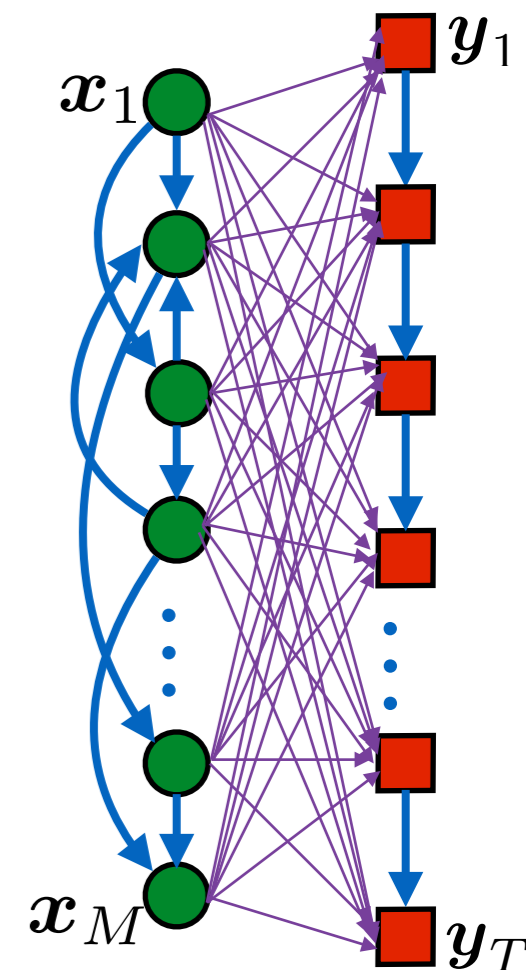
$$\begin{aligned} \max_{\{z_{i,t}\}} \quad & \sum_{t=1}^T \sum_{i=1}^M -z_{i,t} d_{i,t} - \lambda \sum_{i=1}^M \|[z_{i,1} \cdots z_{i,T}]\|_{\infty} \\ & + \beta \left(\sum_{i=1}^M z_{i,1} \log p_1(\mathbf{x}_i) + \sum_{t=2}^T \sum_{i,i'=1}^M z_{i,t-1} z_{i',t} \log p(\mathbf{x}_{i'} | \mathbf{x}_i) \right) \\ \text{s. t.} \quad & z_{i,t} \in \{0, 1\}, \quad \sum_{i=1}^M z_{i,t} = 1, \quad \forall i, t. \end{aligned}$$



Sequential facility location

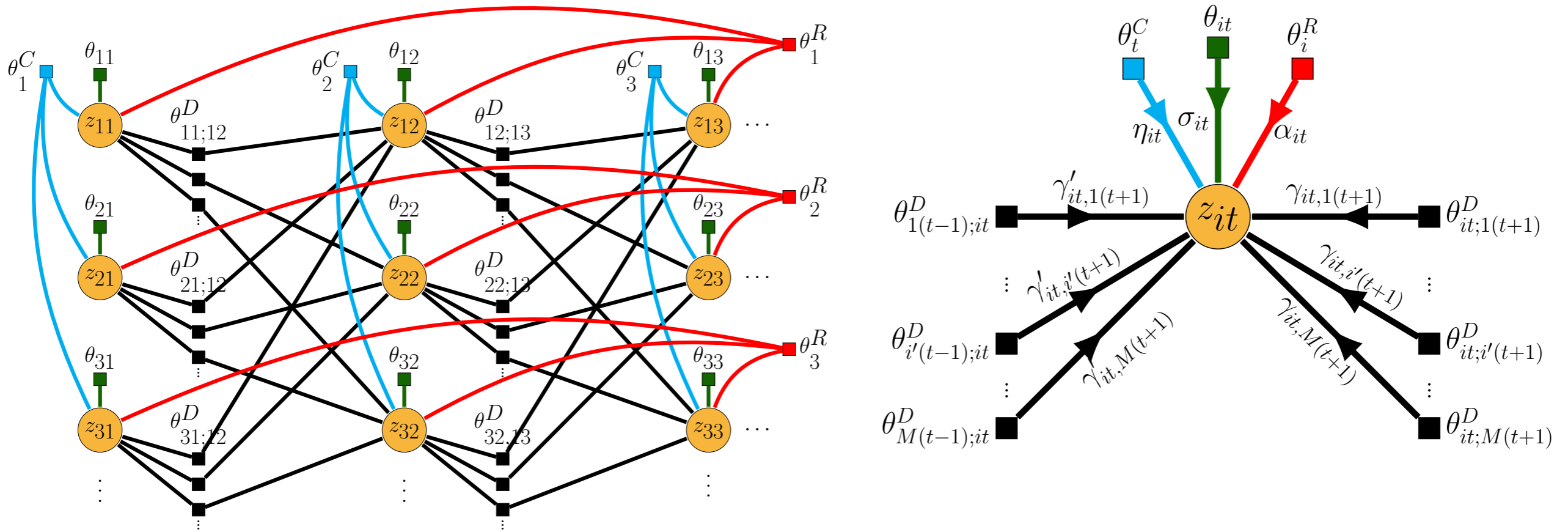
- **SeqFL**: Solve optimization on assignment variables $\{z_{i,t}\}$

$$\begin{aligned} \max_{\{z_{i,t}\}} \quad & \sum_{t=1}^T \sum_{i=1}^M -z_{i,t} d_{i,t} - \lambda \sum_{i=1}^M \|[z_{i,1} \cdots z_{i,T}]\|_{\infty} \\ & + \beta \left(\sum_{i=1}^M z_{i,1} \log p_1(\mathbf{x}_i) + \sum_{t=2}^T \sum_{i,i'=1}^M z_{i,t-1} z_{i',t} \log p(\mathbf{x}_{i'} | \mathbf{x}_i) \right) \\ \text{s. t.} \quad & \underline{z_{i,t} \in \{0, 1\}}, \quad \sum_{i=1}^M z_{i,t} = 1, \quad \forall i, t. \end{aligned}$$



- **Non-convex**: does not make sense (here) to relax!
 - Solve via **Max-Sum Message Passing**

SeqFL via message passing



Encoding factor:

$$\theta_{it}(z_{i,t}) = -\bar{d}_{i,t} z_{i,t}$$

Cardinality factor:

$$\theta_i^R(\mathbf{z}_{i,:}) = \begin{cases} -\lambda, & \|\mathbf{z}_{i,:}\|_\infty > 0 \\ 0, & \text{otherwise} \end{cases}$$

Dynamic factor:

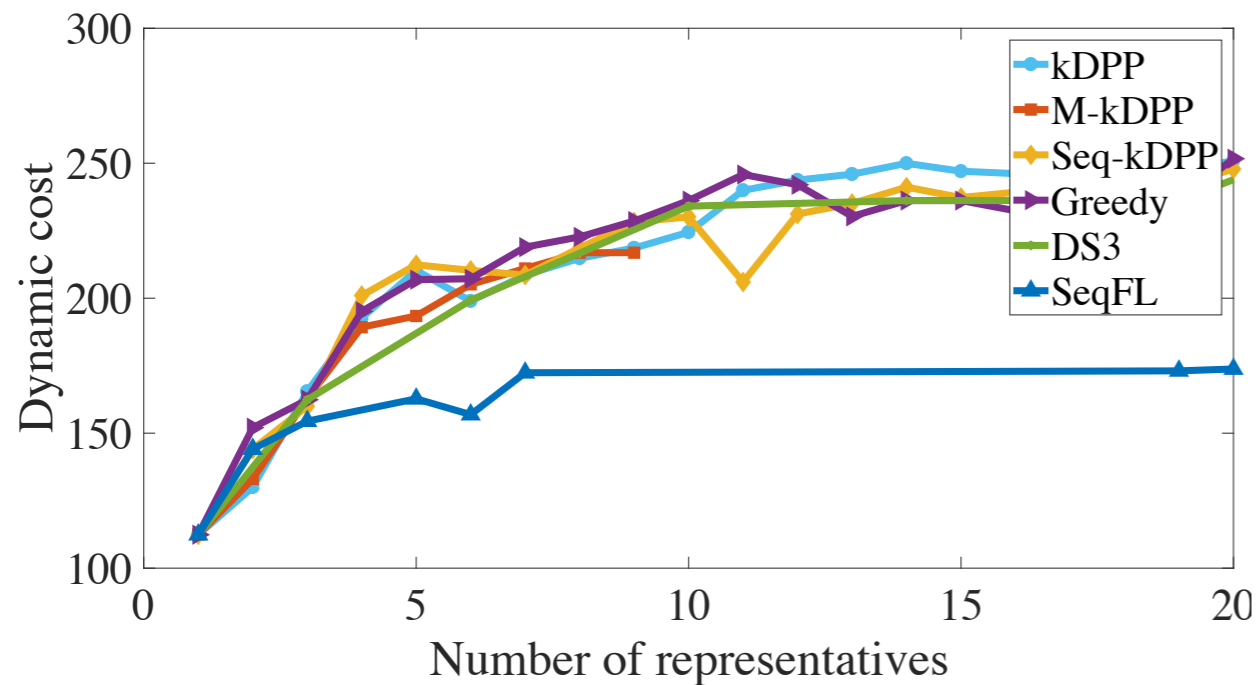
$$\theta_{it-1;i',t}^D(z_{i,t-1}, z_{i',t}) = \log p(\mathbf{x}_{i'} | \mathbf{x}_i) z_{i,t-1} z_{i',t}$$

Constraint factor:

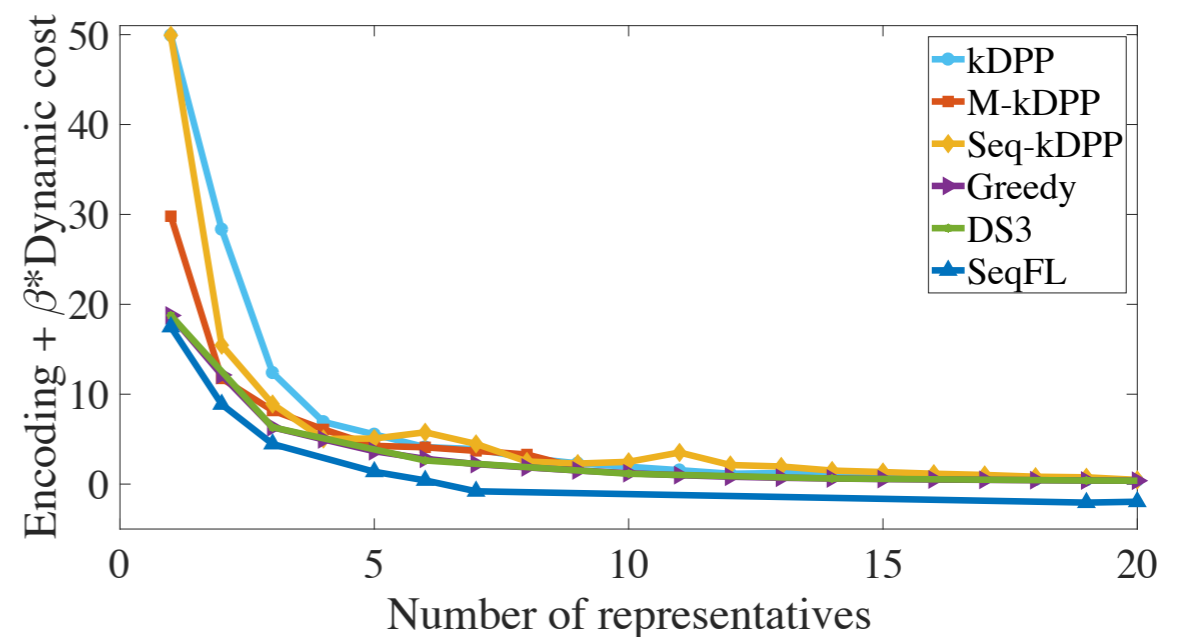
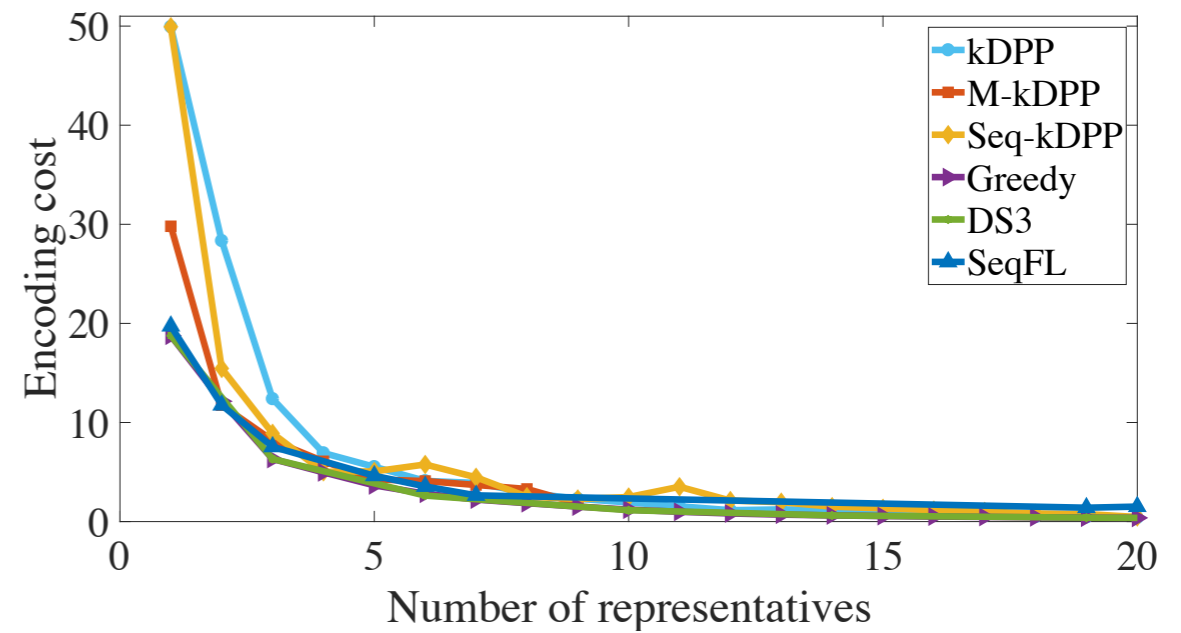
$$\theta_t^C(\mathbf{z}_{:,t}) = \begin{cases} 0, & \sum_{i=1}^M z_{i,t} = 1 \\ -\infty, & \text{otherwise.} \end{cases}$$

Synthetic experiments

- Data from HMM with $M = 50$ states and $T = 100$



Sequence of reps in SeqFL:
compatible with dynamic model



Procedure learning

- We learn **how to do** tasks using **instructional videos**

A screenshot of the YouTube search interface. The search bar contains the text "how to". Below the search bar, it displays "About 547,000,000 results".

- How to *perform CPR*
- How to *change a flat tire*
- How to *make a carrot cake*
- How to *assemble a bike*
- How to *install chrome cast*

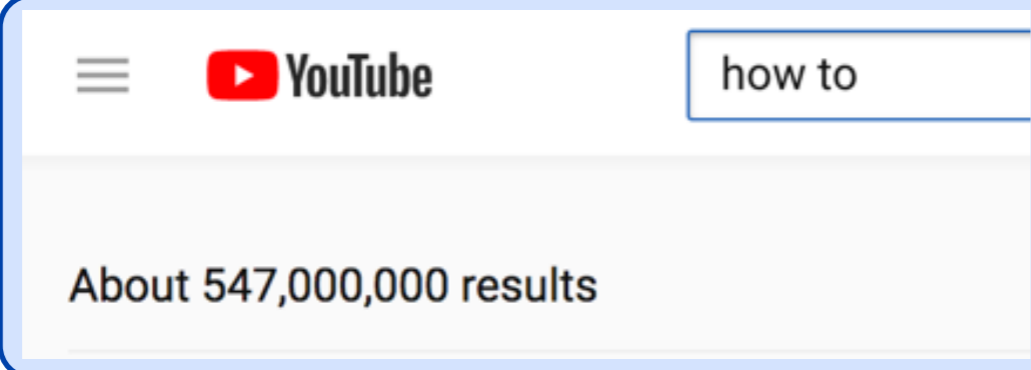
Large
videos

A screenshot of the YouTube search results for "how to perform CPR". The search bar contains the text "how to perform CPR". Below the search bar, it displays "About 97,800 results". The results list includes:

- How to Perform CPR video**
CPRCertified.com • 2.2M views • 3 years ago
This training video shows you how to perform CPR on an adult. Learn what lifesaving measures you can do to save someone's life
4:59
- How To Do CPR - Animated Video**
HealthSketch • 113K views • 1 year ago
A simple animated video explaining how to deliver CPR to people who have collapsed and are not breathing, using the steps "DR'S"
CC
4:03
- How to Give CPR | First Aid Training**
Howcast • 347K views • 5 years ago

Procedure learning

- We learn **how to do** tasks using **instructional videos**



- How to *perform CPR*
- How to *change a flat tire*
- How to *make a carrot cake*
- How to *assemble a bike*
- How to *install chrome cast*

- Goal: **Learn grammar** of complex tasks using video data

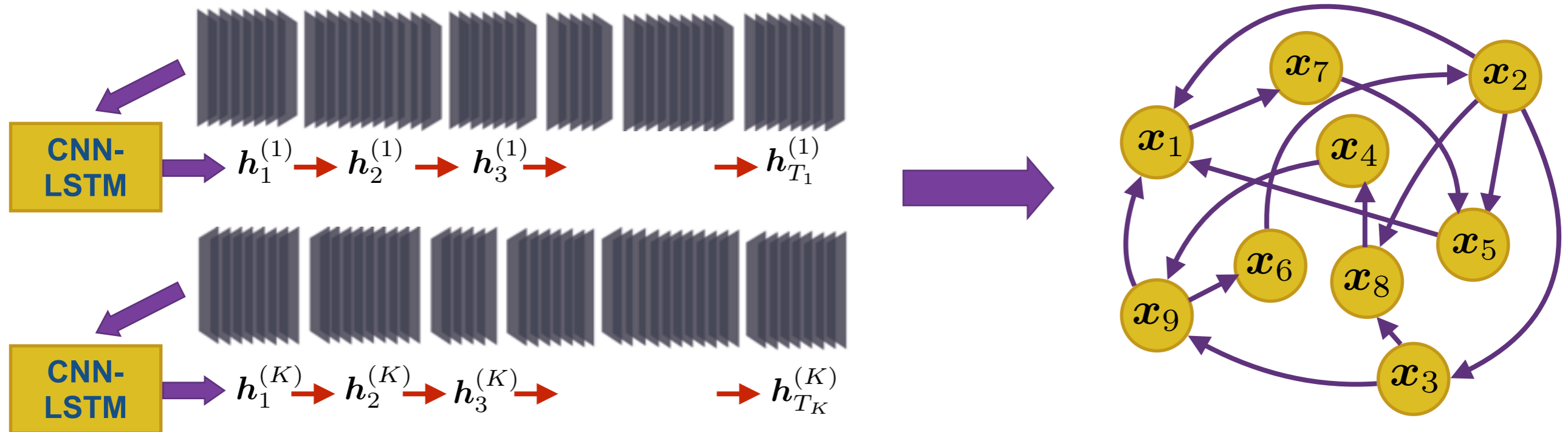


procedure:



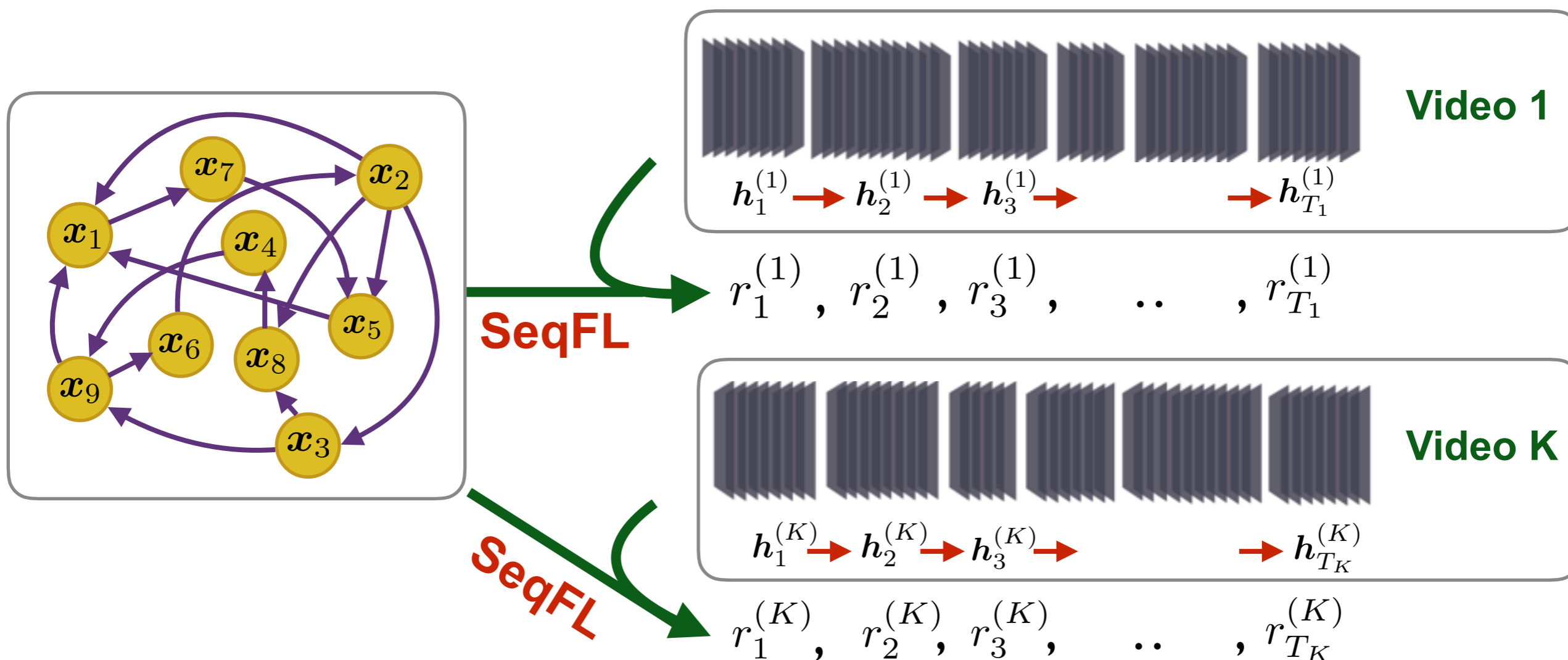
Procedure learning via SeqFL

- 1) Extract **superframes** (Gygli et al '14) from each video
- 2) Extract **deep features** for each superframe
- 3) Build a dynamic model: **1st order HMM**



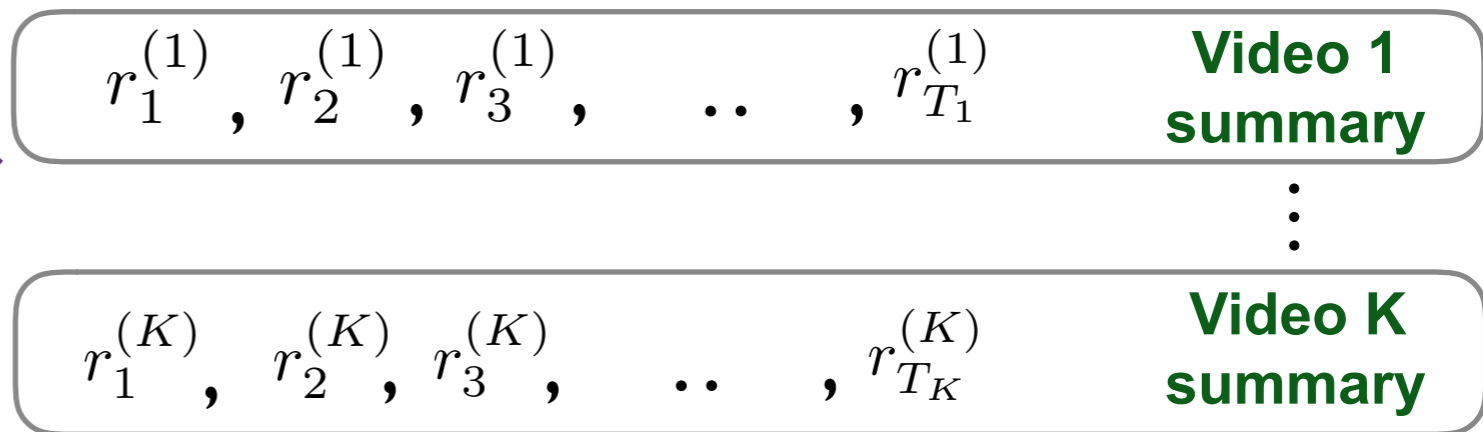
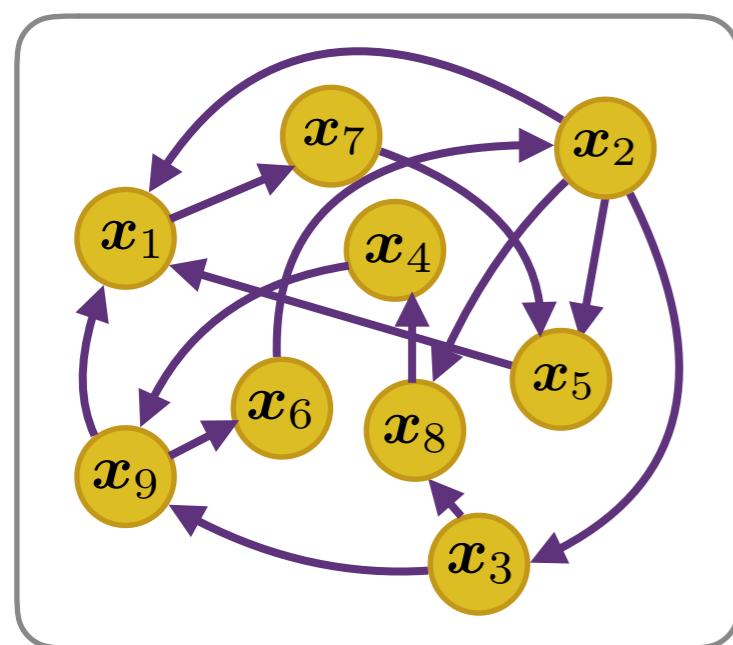
Procedure learning via SeqFL

- 1) Extract **superframes** (Gygli et al '14) from each video
- 2) Extract **deep features** for each superframe
- 3) Build a dynamic model: **1st order HMM**
- 4) **SeqFL** (X = HMM, Y = test videos)

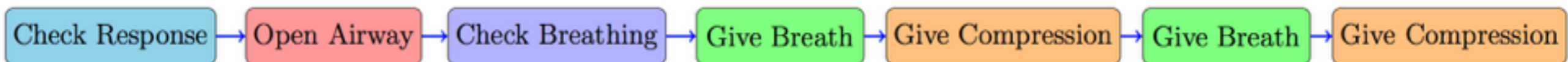


Procedure learning via SeqFL

- 1) Extract **superframes** (Gygli et al '14) from each video
- 2) Extract **deep features** for each superframe
- 3) Build a dynamic model: **1st order HMM**
- 4) **SeqFL** ($X = \text{HMM}$, $Y = \text{test videos}$)
- 5) **Align sequences** of representatives
- 6) **Map** representatives to **actions**



$x_5 \rightarrow x_2 \rightarrow x_9 \rightarrow x_3 \rightarrow x_8 \rightarrow x_3 \rightarrow x_8$



Experimental results

- **Instructional video dataset** (Alayrac et al CVPR '16)
 - **5 tasks, 30 videos/task**, labels of key steps given

Task		kDPP	M-kDPP	Seq-kDPP	DS3	SeqFL
Change tire	(P, R)	(0.56, 0.50)	(0.55, 0.60)	(0.44, 0.40)	(0.56, 0.50)	(0.60, 0.60)
	F-score	0.53	0.57	0.42	0.53	0.60
Make coffee	(P, R)	(0.38, 0.33)	(0.50, 0.44)	(0.63, 0.56)	(0.50, 0.56)	(0.50, 0.56)
	F-score	0.35	0.47	0.59	0.53	0.53
CPR	(P, R)	(0.71, 0.71)	(0.71, 0.71)	(0.71, 0.71)	(0.71, 0.71)	(0.83, 0.71)
	F-score	0.71	0.71	0.71	0.71	0.77
Jump car	(P, R)	(0.50, 0.50)	(0.56, 0.50)	(0.56, 0.50)	(0.50, 0.50)	(0.60, 0.60)
	F-score	0.50	0.53	0.53	0.50	0.60
Repot plant	(P, R)	(0.57, 0.67)	(0.60, 0.50)	(0.57, 0.67)	(0.57, 0.67)	(0.80, 0.67)
	F-score	0.62	0.55	0.62	0.62	0.73
All tasks	(P, R)	(0.54, 0.54)	(0.58, 0.55)	(0.58, 0.57)	(0.57, 0.59)	(0.67, 0.63)
	F-score	0.54	0.57	0.57	0.58	0.65

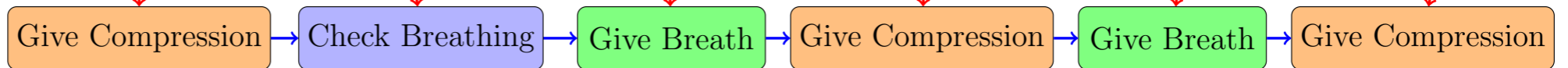

No dynamics
with dynamics

Experimental results

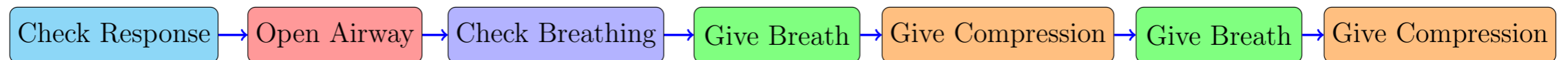
- How to **'perform CPR'**



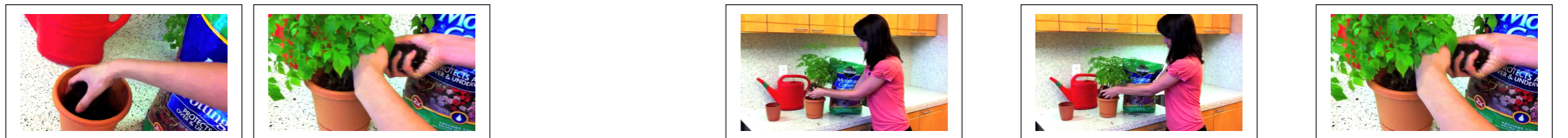
SeqFL



Ground Truth



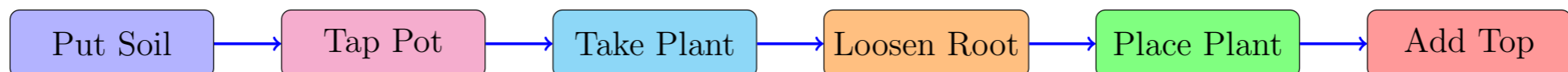
- How to **'repot a plant'**



SeqFL

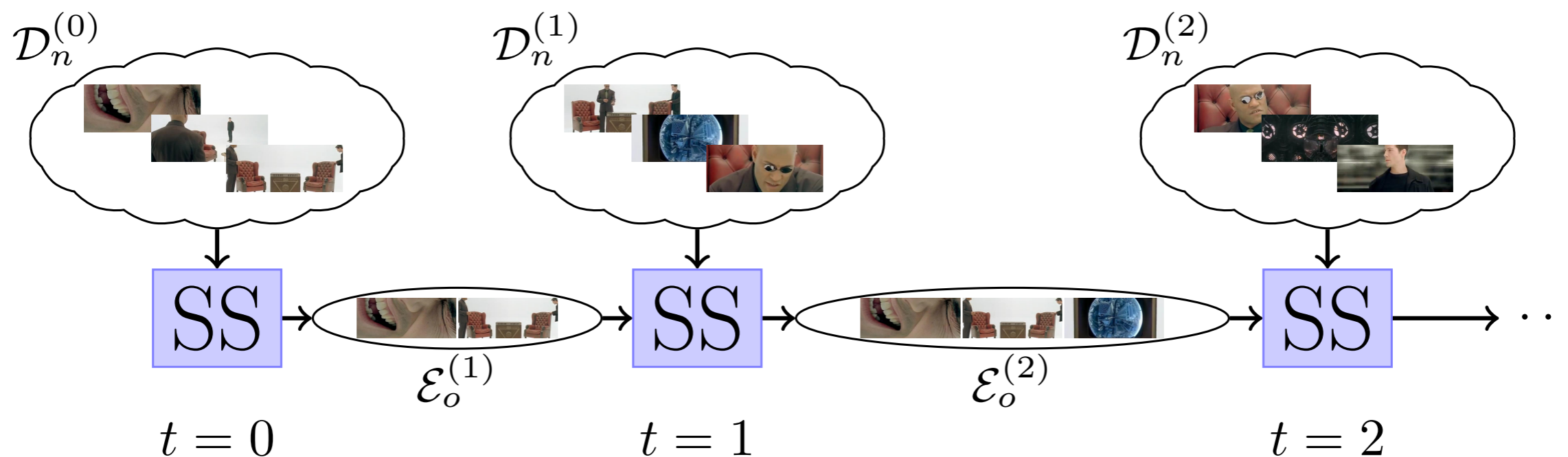


Ground Truth



Online summarization

- **Sequential data**: grow **incrementally**, collecting all and running SS **infeasible**
- Approach: **incremental subset selection**



- Instead of running SS on $\bigcup_t \mathcal{D}_n^{(t)}$, run SS on $\mathcal{E}_o^{(t)} \cup \mathcal{D}_n^{(t)}$ at t
- Proved exactness of the online SS

Summary and ongoing work

- Important to have summarization close to **global optimum**
 - Developed **efficient sparse optimization** for **facility location**
 - Showed **exactness** under appropriate condition
 - Important to handle **sequential data**: **compatible seq of reps**
 - Developed **SeqFL** for sequential summarization
 - Addressed **procedure learning** from **instructional videos**
-

Thanks!

Codes: <http://www.ccs.neu.edu/home/eelhami/codes.htm>