

# Distributed Adaptive Sampling for Kernel Matrix Approximation

Daniele Calandriello, Alessandro Lazaric, Michal Valko



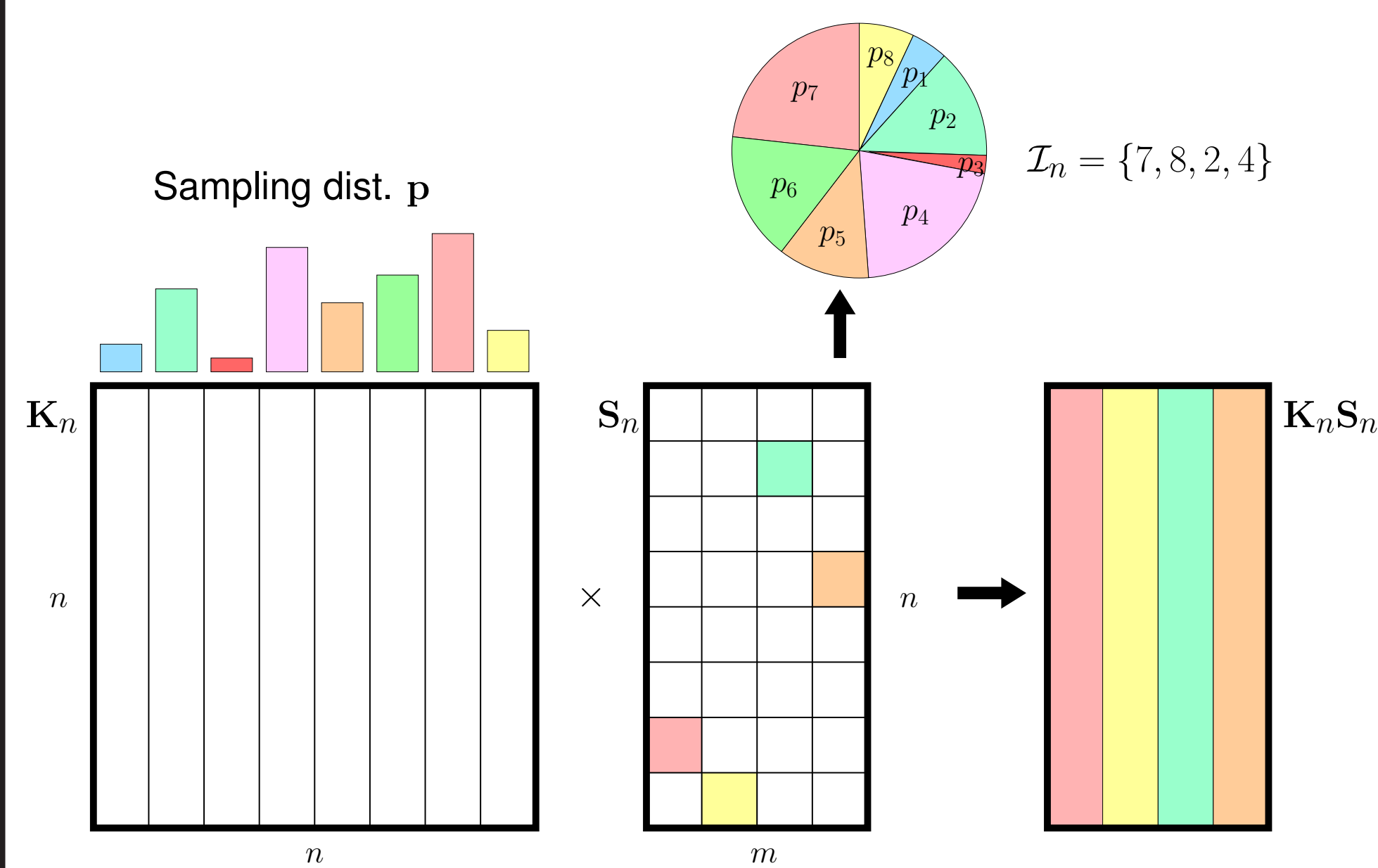
## Motivation

- Kernel methods are *versatile* and *accurate*
- Strong generalization guarantees but *poor scalability*
- $\mathcal{O}(n^3)$  time  $\mathcal{O}(n^2)$  space ( $n$  number of samples)
- Current limitation: Many approximate schemes are either *not scalable* or *not accurate*
- ⇒ We propose a **parallel distributed incremental approximation** scheme for kernel methods with **complexity and error guarantees adaptive to the dataset and kernel structure**
- ⇒ Runs in a **single pass** over the dataset, and can **update its solution** if **new data** arrives
- ⇒ On a single machine, computes a solution in **subquadratic**  $\mathcal{O}(n)$  time, avoids constructing the whole  $K_n$
- ⇒ On multiple machines, computes a solution in **logarithmic**  $\mathcal{O}(\log(n))$  time, without increase in total work
- ⇒ **Black-box applicability** to many downstream tasks

## Nyström Approximation

### Subsampling

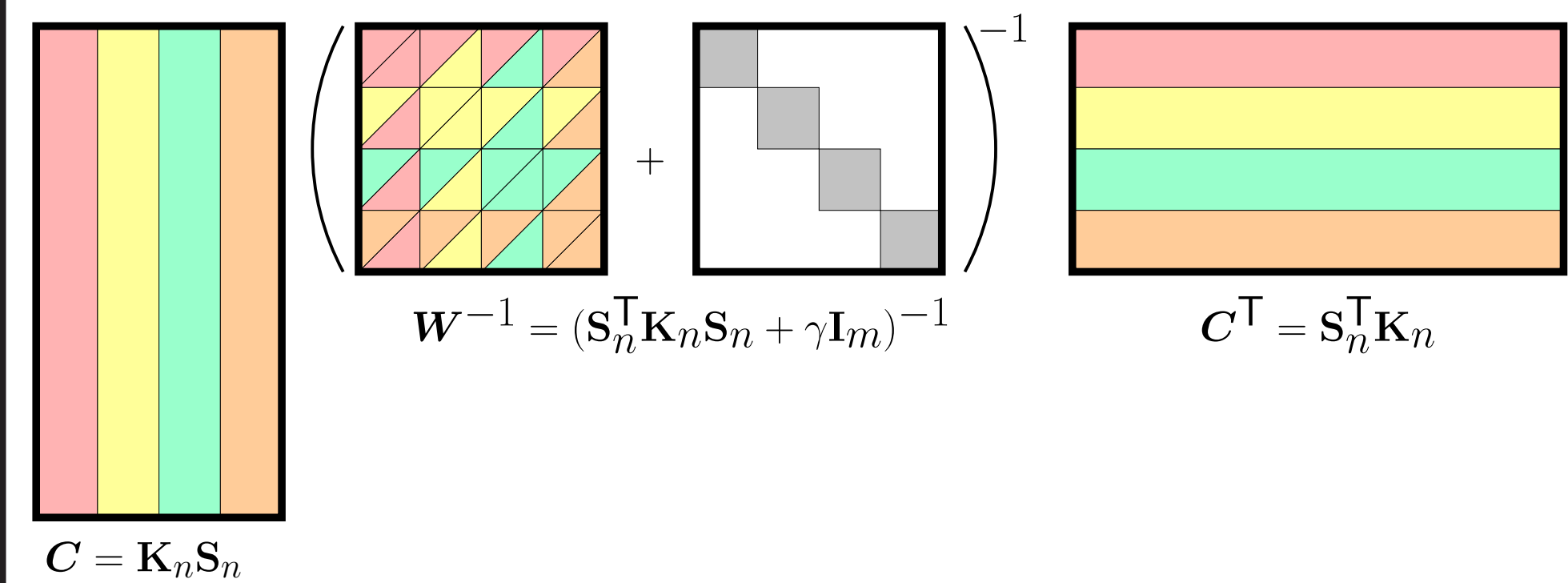
- Select a subset (dictionary)  $\mathcal{I}_n$  of  $m$  representative samples
- Constructs a sparse matrix  $S_n$  to select and reweight the columns associated with the points in  $\mathcal{I}_n$



### Low-Rank Approximation

- Compute approximate, low-rank matrix  $\tilde{K}_n = CW^{-1}C^T$  as

$$\tilde{K}_n = CW^{-1}C^T = K_n S_n (S_n^T K_n S_n + \gamma I_m)^{-1} S_n^T K_n$$



### Efficient Solution

- Compute approximate solution Kernel Ridge Regression

$$\tilde{w}_n = (\tilde{K}_n + \mu I)^{-1} y_n = \frac{1}{\mu} (y_n - C(C^T C + \mu W)^{-1} C^T y_n)$$

### Kernel K-Means

$$\min_A \text{Tr}(K_n - AA^T K_n AA^T) \sim \min_A \text{Tr}(\tilde{K}_n - AA^T \tilde{K}_n AA^T)$$

### Kernel PCA

$$\min_Z \|K_n - ZZ^T K_n\|_F \sim \min_Z \|\tilde{K}_n - ZZ^T \tilde{K}_n\|_F$$

Also Kernel CCA, Kernel [Your downstream problem here]

### Scalability

now depends on  $m$

Space:  $\mathcal{O}(n^2) \Rightarrow \mathcal{O}(nm)$ , Time:  $\mathcal{O}(n^3) \Rightarrow \mathcal{O}(nm^2 + m^3)$

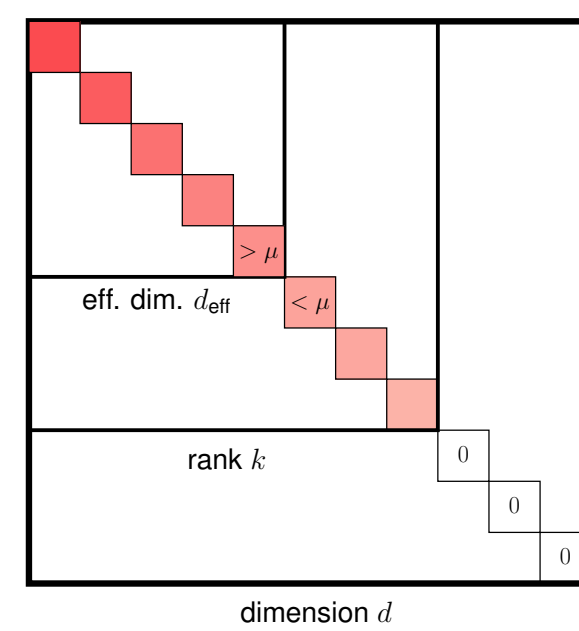
### Problems:

- ? How to choose the sampling distribution?
- ? How to choose  $m$ ?

## References

- [Alaoui and Mahoney (2015)] A. El Alaoui and M. W. Mahoney. Fast randomized kernel methods with statistical guarantees. In *NIPS*, 2015.
- [Bach (2013)] F. Bach. Sharp analysis of low-rank kernel matrix approximations. In *COLT*, 2013.
- [Calandriello et al. (2016)] D. Calandriello, A. Lazaric, and M. Valko. Analysis of Nyström method with sequential ridge leverage scores. In *UAI*, 2016.
- [Rudi et al. (2015)] A. Rudi, R. Camoriano, and L. Rosasco. Less is more: Nyström computational regularization. In *NIPS*, 2015.
- [Musco and Musco (2016)] C. Musco and C. Musco. Provably useful kernel matrix approximation in linear time. In arXiv, 2016.

## Kernel Ridge Leverage Scores (RLS) Sampling



Definition 1. Given a kernel matrix  $K_n \in \mathbb{R}^{n \times n}$ , define

$$\gamma\text{-ridge leverage score} \quad \tau_{n,i}(\gamma) = e_{n,i}^T K_n^{-1} (K_n + \gamma I_n)^{-1} e_{n,i} = \phi(x_i)^T (\phi(X_n) \phi(X_n)^T + \gamma I)^{-1} \phi(x_i) \quad (1)$$

$$\text{effective dimension} \quad d_{\text{eff}}(\gamma)_n = \sum_{i=1}^n \tau_{n,i}(\gamma) = \text{Tr}(K_n (K_n + \gamma I_n)^{-1}) \quad (2)$$

Proposition 1 (Alaoui, Mahoney, 2015). Let  $\epsilon$  be the accuracy,  $\delta$  the confidence. If the regularized Nyström approximation  $\tilde{K}_n$  is computed using a sampling distribution proportional to  $\tau_{n,i}$ , and at least

$$m \geq \left( \frac{2 d_{\text{eff}}(\gamma)_n}{\epsilon^2} \right) \log \left( \frac{n}{\delta} \right)$$

columns, then with probability  $1 - \delta$ ,  $0 \preceq K_n - \tilde{K}_n \preceq \frac{\gamma}{1-\epsilon} I_n$ .

Intuitively:  $\tau_{n,i}$  sensitivity of prediction on point  $x_i$   
 $\Rightarrow \hat{y}_{n,i} = e_i^T (K_n \hat{w}_n) = e_i^T K_n (K_n + \mu I)^{-1} y_n$

Pros: +  $m$  scales with the effective dimension

Cons: - computing  $\tau_{n,i}(\mu)$  is as difficult as solving the original problem  
 - the probabilities need be recomputed at any new sample (=multipass)

## SQUEAK

Lemma 1. Assume that the dictionary  $\mathcal{I}_{t-1}$  induces a  $\gamma$ -approx.  $\tilde{K}_{t-1}$ , and let  $\tilde{S}_t$  be constructed by adding  $\bar{q}$  copies of  $(\bar{q})^{-1/2} e_{t,t}$  to the selection matrix. Then, denoting  $\alpha = (1 + \epsilon)/(1 - \epsilon)$ , for all  $i$  such that  $i \in \mathcal{I}_{t-1} \cup \{t\}$ ,

$$\tilde{\tau}_{t,i} = \frac{1 + \epsilon}{\alpha \gamma} \left( k_{i,i} - k_{t,i} \tilde{S}^T (\tilde{S}^T K_t \tilde{S} + \gamma I)^{-1} \tilde{S}^T k_{t,i} \right), \quad (3)$$

is an  $\alpha$ -approximation of the RLS  $\tau_{t,i}$ , that is  $\tau_{t,i}(\gamma)/\alpha \leq \tilde{\tau}_{t,i} \leq \tau_{t,i}(\gamma)$ .

### SQUEAK

Input:  $\mathcal{D}$ , regularization  $\gamma, \bar{q}, \epsilon$ , Output:  $\mathcal{I}_n$

- Initialize  $\mathcal{I}_0$  as empty,  $\tilde{p}_{1,0} = 1$
- for  $t = 1, \dots, n$  do
- Receive new sample  $x_t$
- Compute  $\alpha$ -app. RLS  $\{\tilde{\tau}_{t,i} : i \in \mathcal{I}_{t-1} \cup \{t\}\}$ , using  $\mathcal{I}_{t-1}$ ,  $x$ , and Eq. 3
- Set  $\tilde{p}_{t,i} = \min\{\tilde{\tau}_{t,i}, \tilde{p}_{t-1,i}\}$
- Initialize  $\mathcal{I}_t = \emptyset$
- for all  $j \in \{1, \dots, t-1\}$  do
- if  $q_{t-1,j} \neq 0$  then
- $q_{t,j} \sim \mathcal{B}(\tilde{p}_{t,j}/\tilde{p}_{t-1,j}, q_{t-1,j})$
- Add  $(j, \phi_j, q_{t,j}, \tilde{p}_{t,j})$  to  $\mathcal{I}_t$ . } SHRINK
- end if
- end for
- $q_{t,t} \sim \mathcal{B}(\tilde{p}_{t,t}, \bar{q})$
- Add  $q_{t,t}$  copies of  $(t, \phi_t, q_{t,t}, \tilde{p}_{t,t})$  to  $\mathcal{I}_t$  } EXPAND
- end for

Theorem 1. Let  $\alpha = (1+\epsilon)/(1-\epsilon)$  and  $\gamma > 1$ . For any  $0 \leq \epsilon \leq 1$ , and  $0 \leq \delta \leq 1$ , if we run SQUEAK with  $\bar{q} = \mathcal{O}(\frac{\alpha}{\epsilon^2} \log(\frac{n}{\delta}))$ , then w.p.  $1 - \delta$ , for all  $t \in [n]$

- $\tilde{K}_t$  computed with  $\mathcal{I}_t$  is a  $\gamma$ -approximation of  $K_t$ .
- $|\mathcal{I}_t| = \sum_i q_{t,i} \leq \mathcal{O}(\bar{q} d_{\text{eff}}(\gamma)_t) \leq \mathcal{O}(\frac{\alpha}{\epsilon^2} d_{\text{eff}}(\gamma)_n \log(\frac{n}{\delta}))$ .

Accuracy and space/time anytime guarantees, matches exact RLS sampling.

Using unnormalized  $\tilde{p}_{t,i}$ , no need for appr.  $d_{\text{eff}}(\gamma)_t$

Only need to compute RLS for points in  $\mathcal{I}_t$ , never recompute after dropping  
 ↳ Never construct the whole  $K_n$ , subquadratic runtime  $\mathcal{O}(n^2 |\mathcal{I}_n|^2) \Rightarrow \mathcal{O}(n |\mathcal{I}_n|^3)$

Store points directly in the dictionary  
 ↳  $\mathcal{O}(d_{\text{eff}}(\gamma)_n^2 + d_{\text{eff}}(\gamma)_n d)$  space constant in  $n$   
 Single pass over the dataset (streaming)

Dictionary changes a lot between iteration, total runtime  $\mathcal{O}(n |\mathcal{I}_n|^3)$

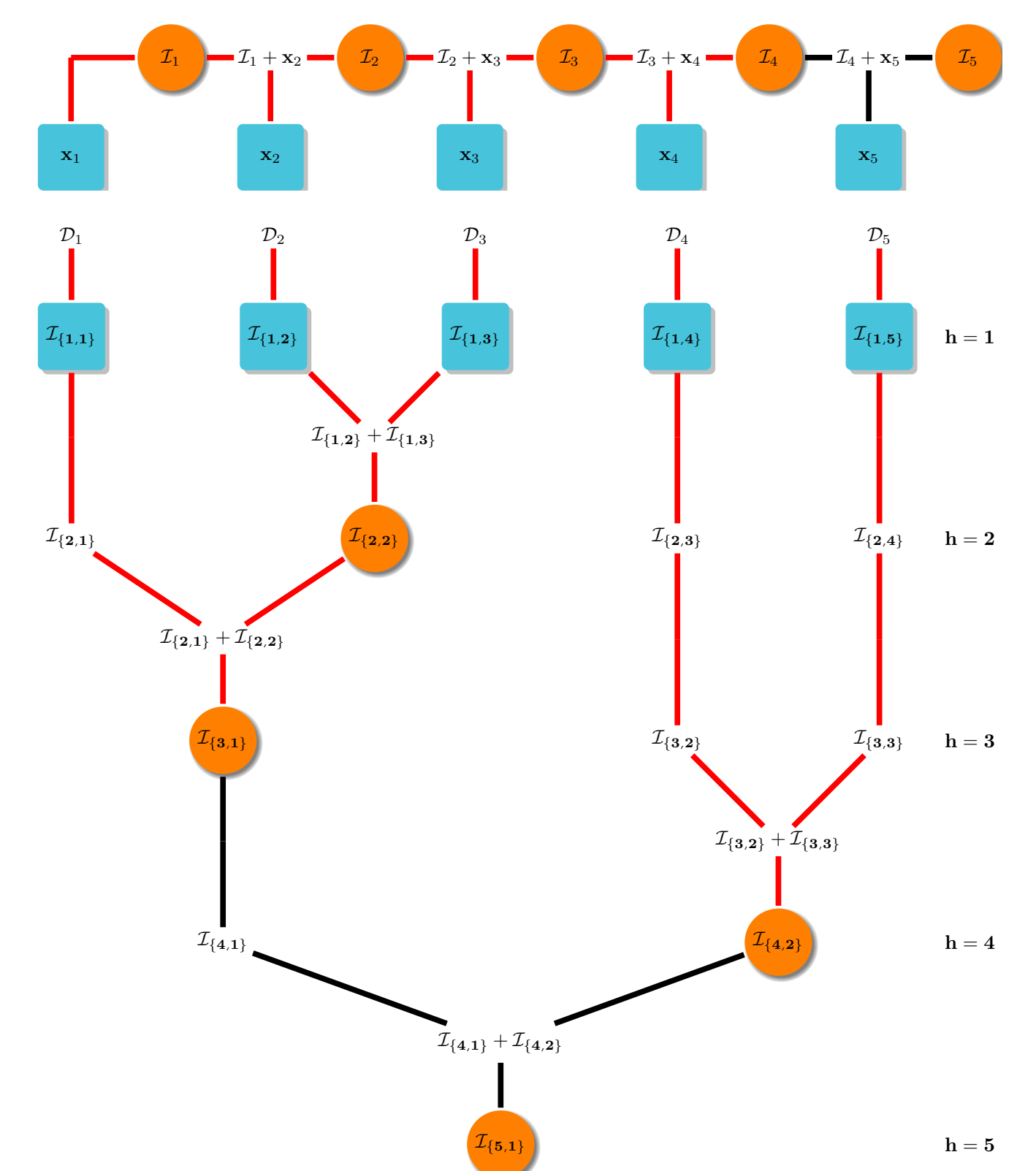
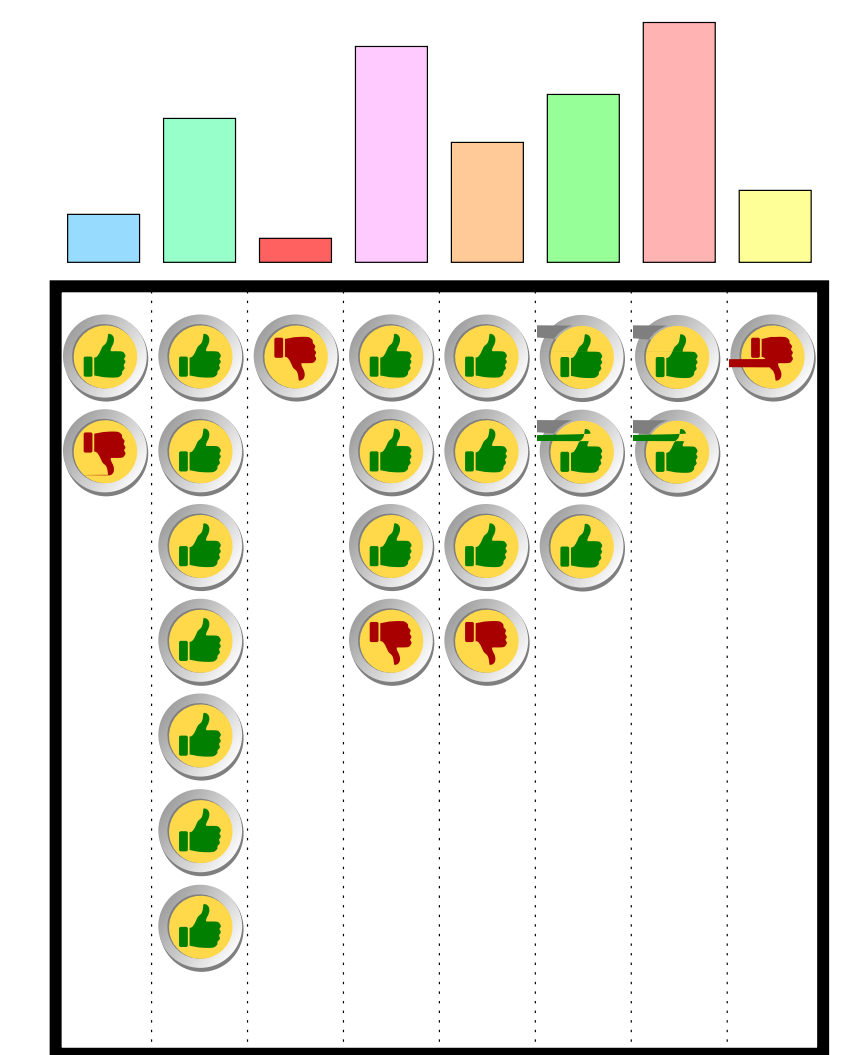
Extend DICT-UPDATE (point + dict.) to DICT-MERGE (dict. + dict.)  
 ↳ Distributed SQUEAK, multiple workers in parallel, without sharing memory  
 Recursive merging to build dictionary,  $\mathcal{O}(\log(n) |\mathcal{I}_n|^3)$  time,  $\mathcal{O}(n |\mathcal{I}_n|^3)$  work

	Time	$ \mathcal{I}_n $	Incr.
EXACT	$n^3$	$n$	-
Bach'13	$\frac{n d_{\text{max},n}^2}{\epsilon} + \frac{d_{\text{max},n}^3}{\epsilon}$	$\frac{d_{\text{max},n}}{\epsilon}$	No
A&M'15	$n( \mathcal{I}_n )^2$	$\left( \frac{\lambda_{\min} + n\mu\epsilon}{\lambda_{\min} - n\mu\epsilon} \right) d_{\text{eff}}(\gamma)_n$	No
INK (C&al'16)	$\frac{\lambda_{\text{max}}^2}{\gamma^2} \frac{n^2 d_{\text{eff}}(\gamma)_n^2}{\epsilon^2}$	$\frac{\lambda_{\text{max}}}{\gamma} \frac{d_{\text{eff}}(\gamma)_n}{\epsilon^2}$	Yes
SQUEAK	$\frac{n^2 d_{\text{eff}}(\gamma)_n^2}{\epsilon^2}$	$\frac{d_{\text{eff}}(\gamma)_n}{\epsilon^2}$	Yes

- $\tilde{\tau}_{t,i} = e_i^T \tilde{K}_t (\tilde{K}_t + \gamma I)^{-1} e_i$  would fail
- Instead, approximate  $\tau_{t,i}$  directly in RKHS  
 $\tilde{\tau}_{t,i} = \phi(x_i)^T (\phi(X_t) \tilde{S} \tilde{S}^T \phi(X_t)^T + \gamma I)^{-1} \phi(x_i)$   
 and then reformulate using kernel trick
- $\tilde{\tau}_{t,i}$  can be computed in  $\mathcal{O}(|\mathcal{I}_t|^2)$  space and  $\mathcal{O}(|\mathcal{I}_t|^3)$  time, independent from  $t$ .
- $\tilde{\tau}_{t,i}$  for samples in  $\mathcal{I}_t$  can be computed using only samples contained in  $\mathcal{I}_t$ .
- The formulation of  $\tilde{\tau}_{t,i}$  is not incremental

Proposition 2. For any kernel matrix  $K_{t-1}$  and its bordering  $K_t$ ,

$$\tau_{t,i} \leq \tau_{t-1,i}, \quad d_{\text{eff}}(\gamma)_t \geq d_{\text{eff}}(\gamma)_{t-1}$$



## Downstream guarantees (Musco & Musco 2016)

RLS sampling preserves well the projection on  $K_n$ 's range  $P = K_n^{1/2} (K_n + \gamma I)^{-1} K_n^{1/2} = \phi(X_n)^T (\phi(X_n) \phi(X_n)^T + \gamma I)^{-1} \phi(X_n)$

Kernel ridge regression:  
 $\hat{y}_{n,i} = e_i^T K_n (K_n + \mu I)^{-1} y_n = e_i^T P y_n$   
 $\tilde{w}_n = (\tilde{K}_n + \gamma I)^{-1} y_n$

Kernel PCA:  
 $K_n = U U^T, P = U (U^T + \gamma I)^{-1} U^T$   
 $\tilde{Z}$  computed using  $\tilde{K}_n = \tilde{U} \tilde{U}^T$

Kernel K-Means:  
 $\tilde{A}$   $\rho$ -optimal cluster assignment for  $\tilde{K}_n$   
 $A^*$  optimal cluster assignment for  $K_n$   
 $\xi = (1 + \epsilon)(1 + \rho)$   
 $\text{Tr}(K_n - \tilde{A} \tilde{A}^T K_n \tilde{A} \tilde{A}^T)$   
 $\leq \xi \text{Tr}(K_n - A^* A^{*T} K_n A^* A^{*T})$

$$R(\tilde{w}_n) \leq \left(1 + \frac{1}{1-\epsilon}\right) R(\hat{w}_n)$$

$$\|K_n - \tilde{Z} \tilde{Z}^T K_n\|_F \leq (1 + 2\epsilon) \|K_n - Z^* Z^{*T} K_n\|_F$$