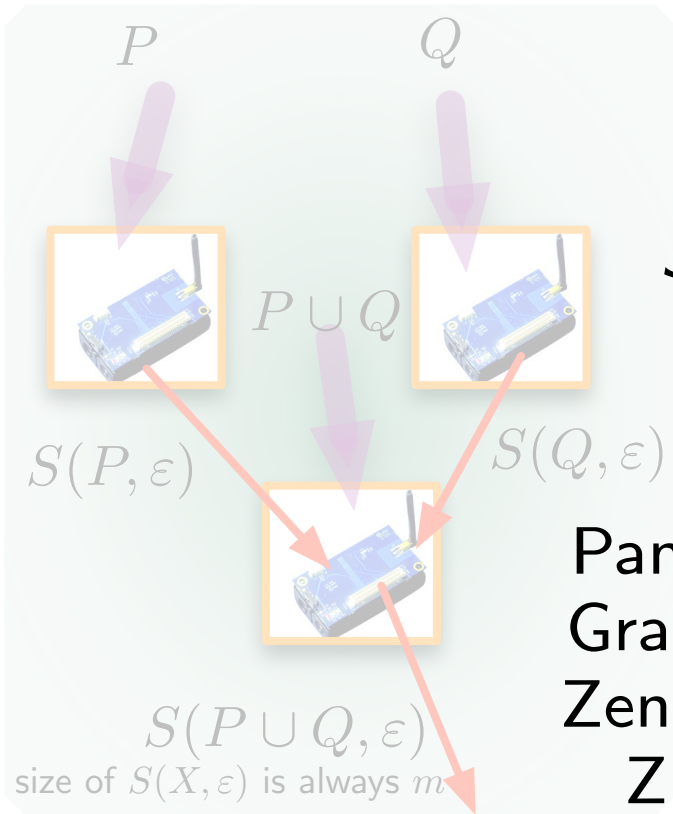
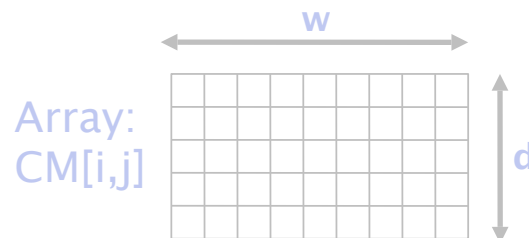
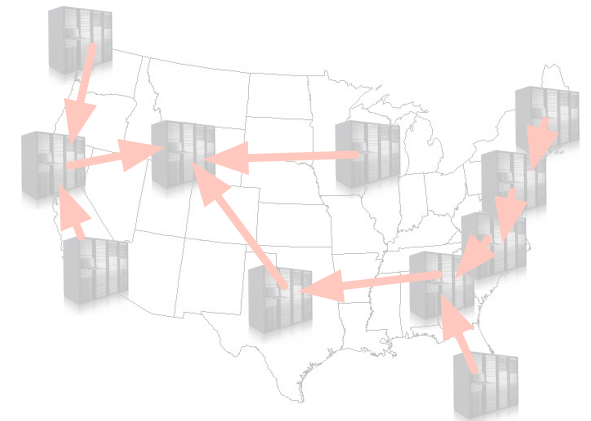
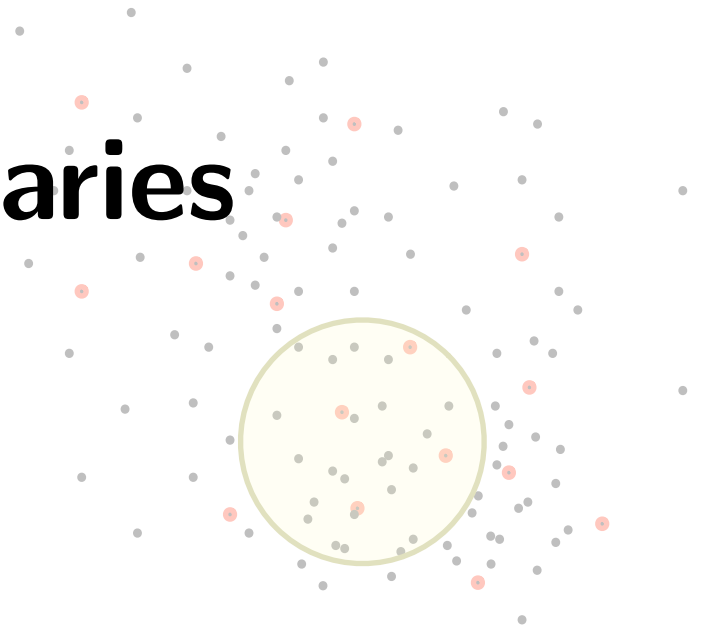


Mergeable Summaries



Jeff M. Phillips
University of Utah

joint with with
Pankaj K. Agarwal (Duke)
Graham Cormode (AT&T)
Zengfeng Huang (HKUST)
Zheiwai Wei (HKUST)
Ke Yi (HKUST)



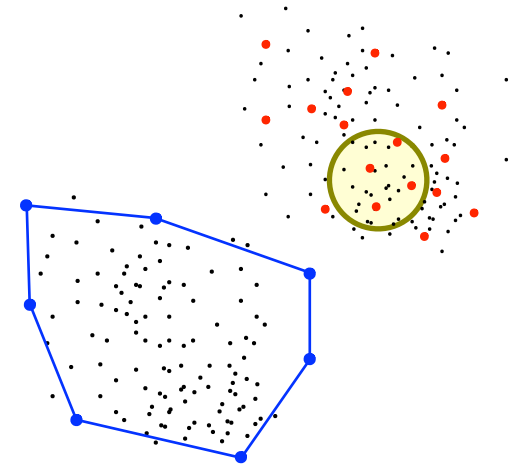
Summaries for MASSIVE Data

Allows approximate computation with guarantees and small space

coreset: small summary, proxy for full data set

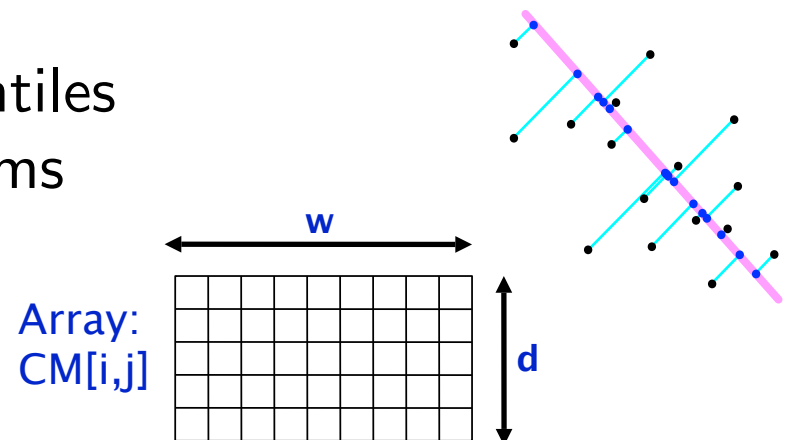
with approx guarantees:

- ε -samples of (P, \mathcal{R}) : approx density
- ε -kernel: approx convex shape



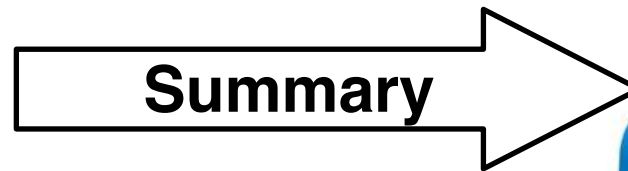
sketch: (random) (linear) combination of full data, recover functions with approx guarantees:

- Euclidean distance: Johnson-Lindenstrauss random projection
- min-count sketch: approx item counts
- Greenwald-Khanna sketch: approx quantiles
- Misra-Gries sketch: approx frequent items



Summaries for MASSIVE Data

Allows approximate computation with guarantees and small space



Massive Distributed Computation

data centers



Massive Distributed Computation

data centers



Massive Distributed Computation

data centers



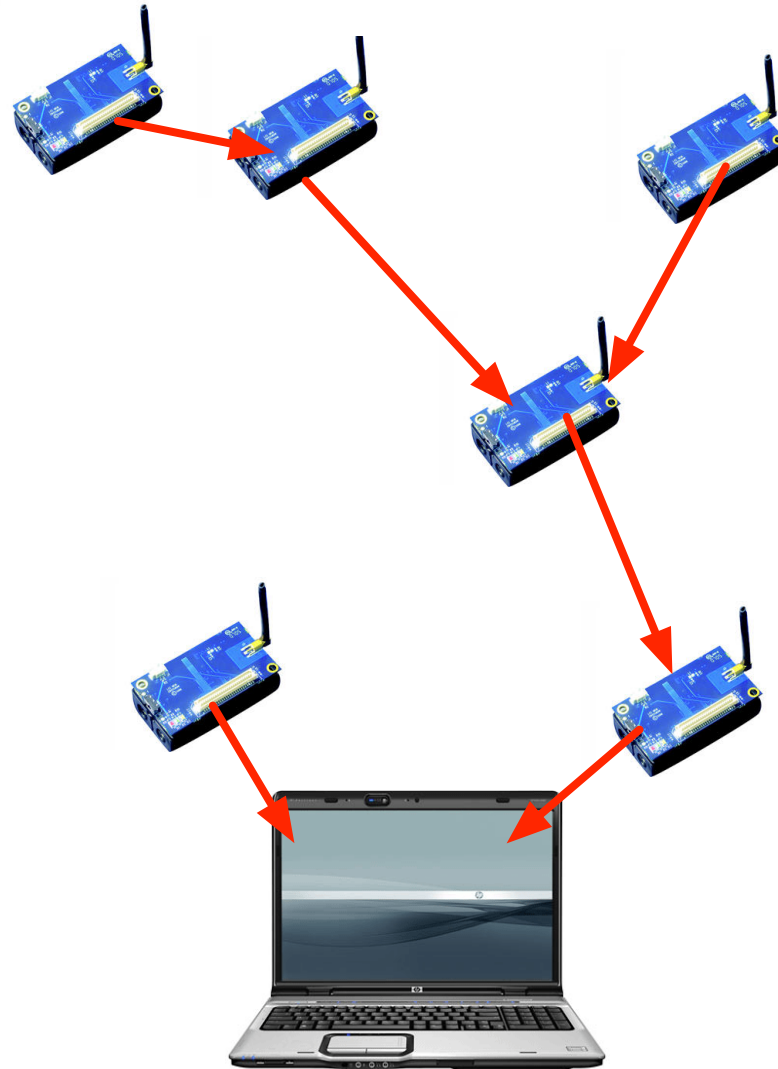
Massive Distributed Computation

data centers



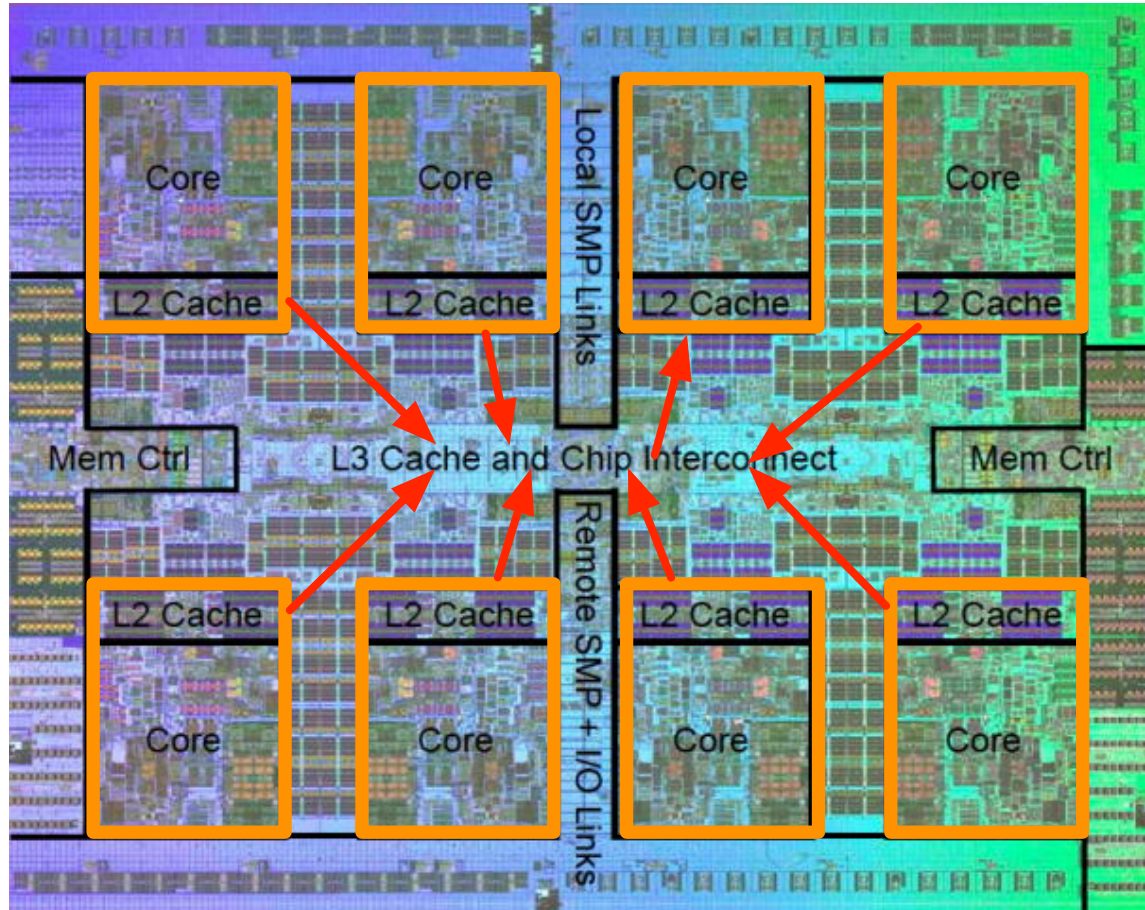
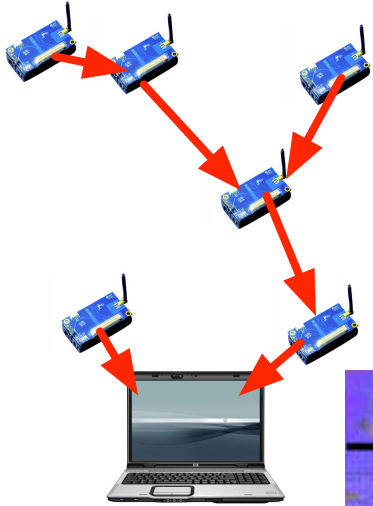
Massive Distributed Computation

data centers
sensor networks



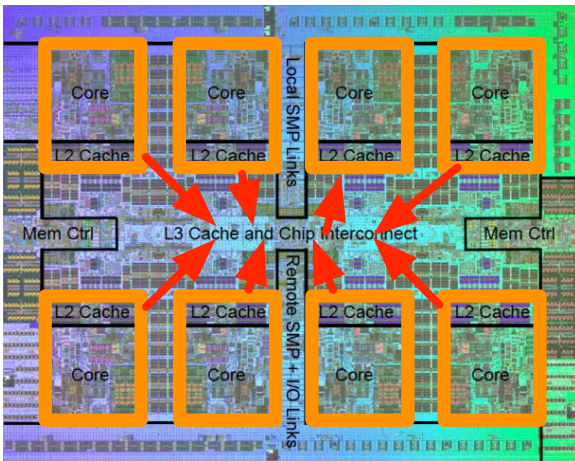
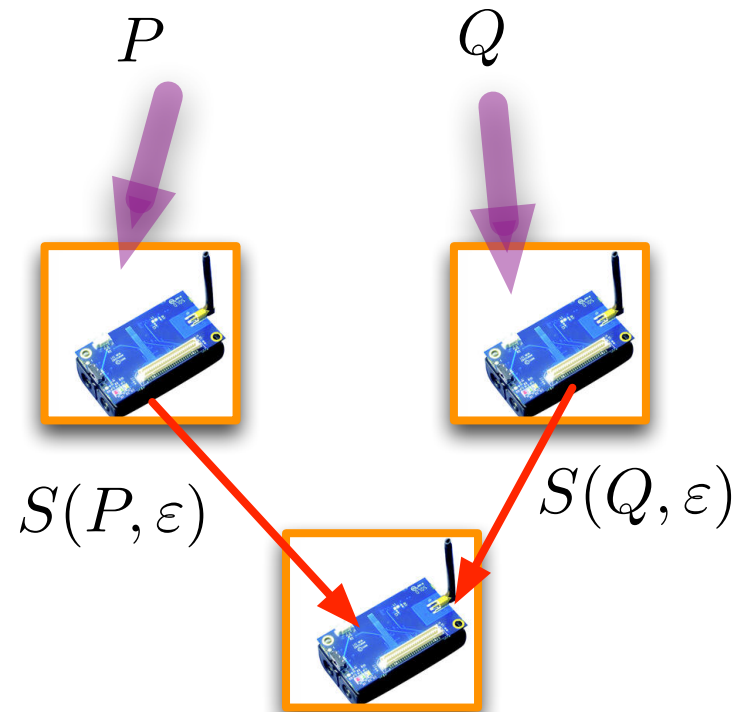
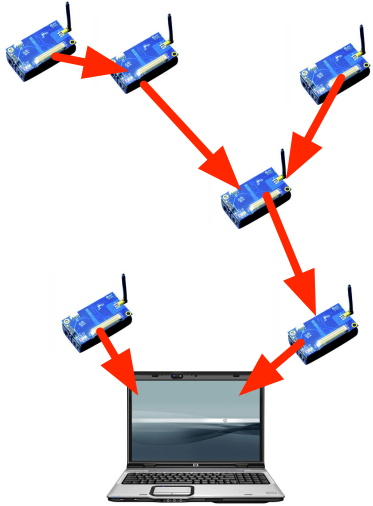
Massive Distributed Computation

data centers
sensor networks
multi-core



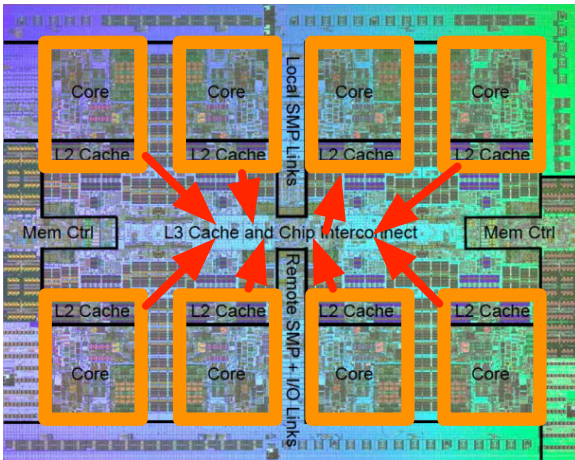
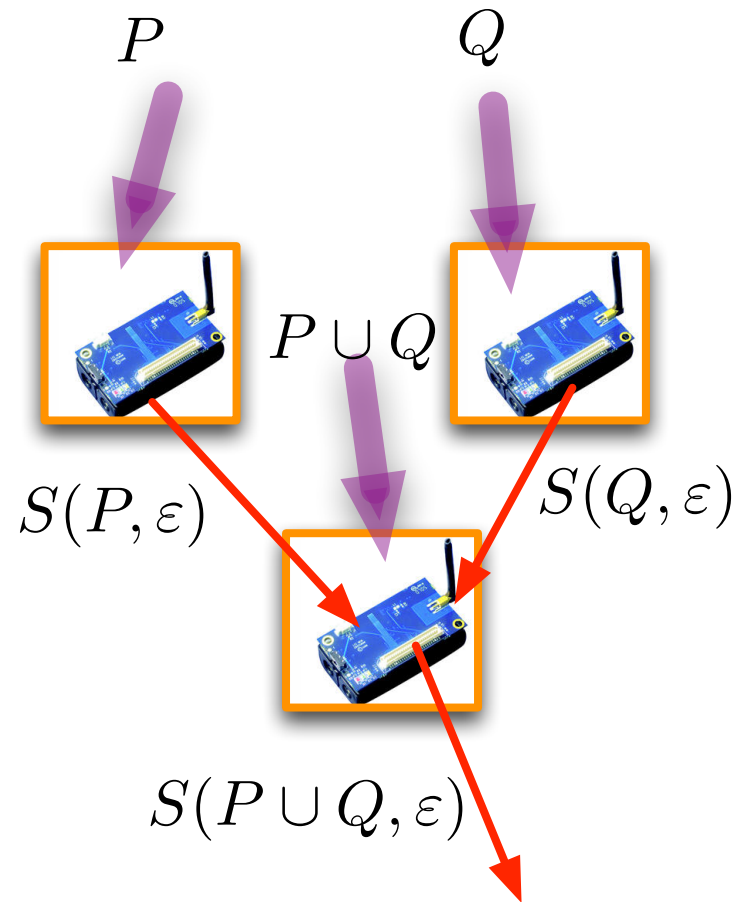
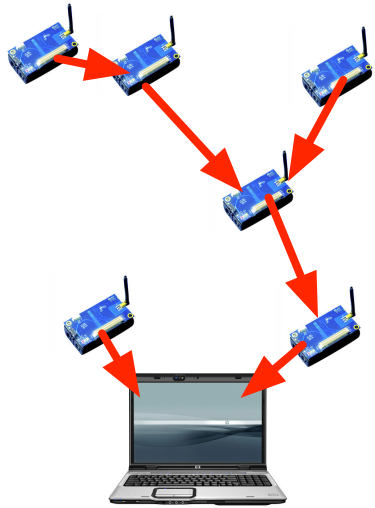
Massive Distributed Computation

data centers
sensor networks
multi-core



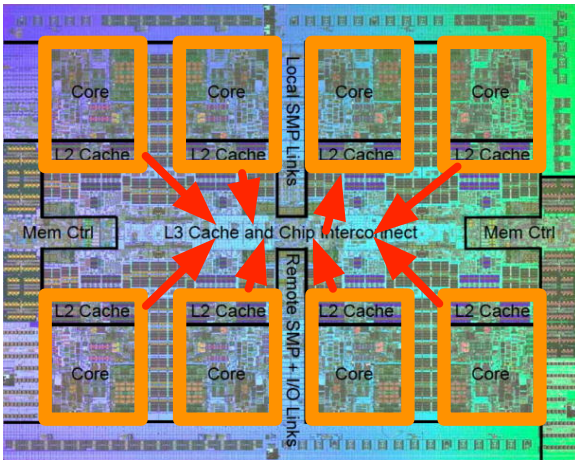
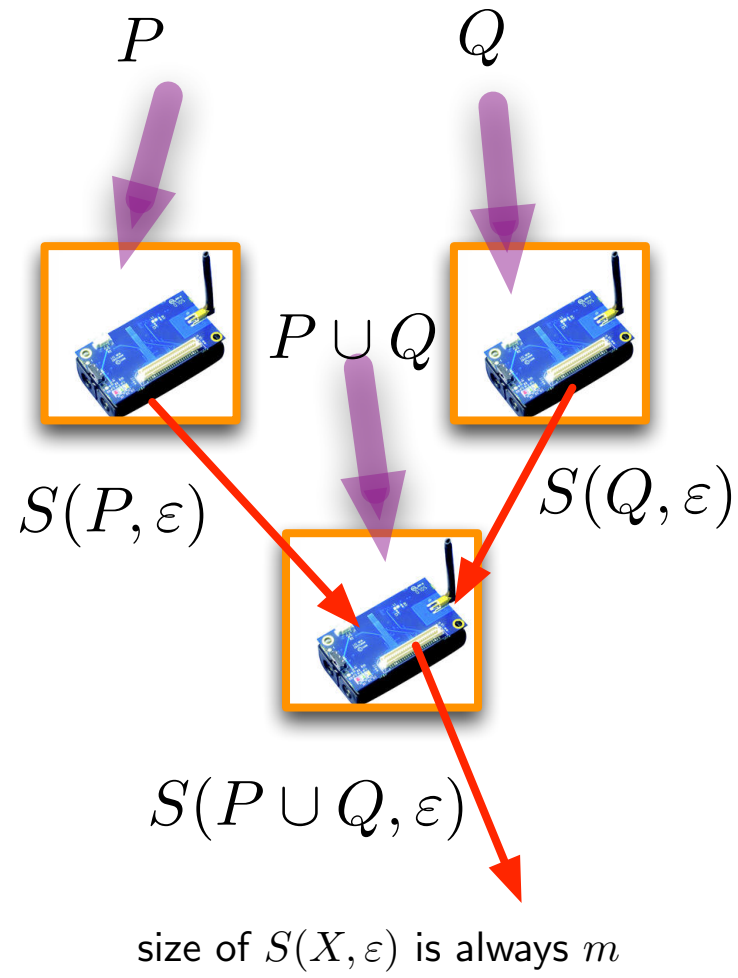
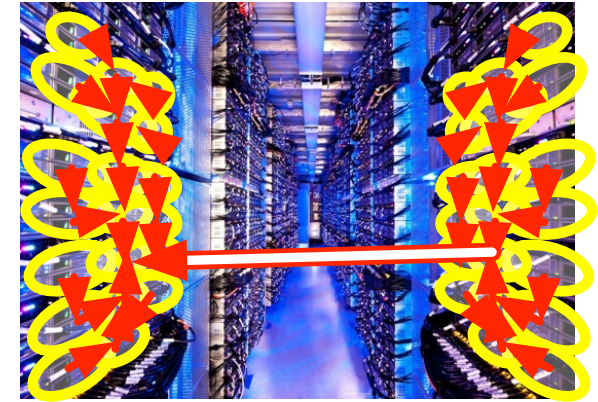
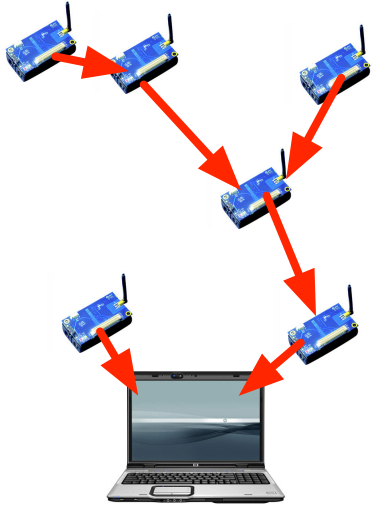
Massive Distributed Computation

data centers
sensor networks
multi-core



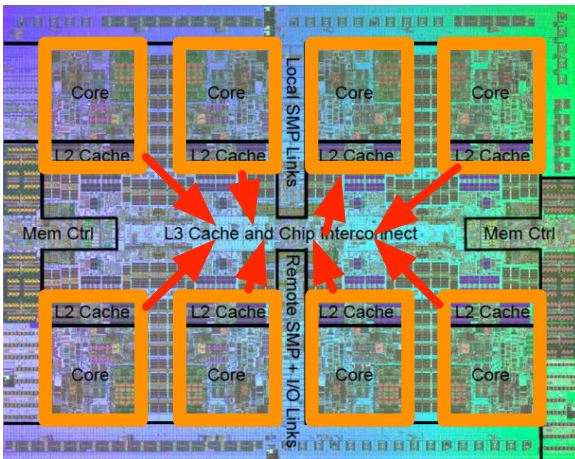
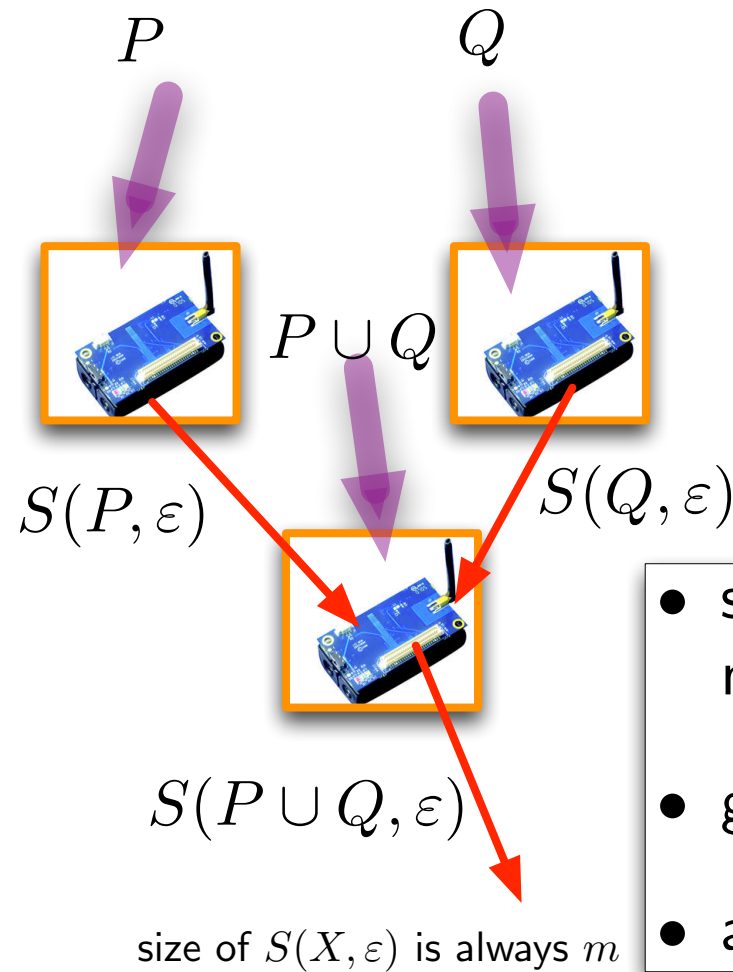
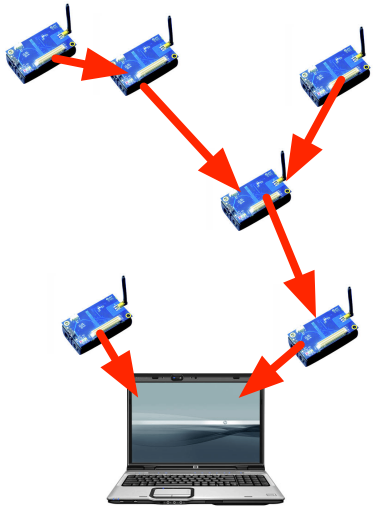
Massive Distributed Computation

data centers
sensor networks
multi-core



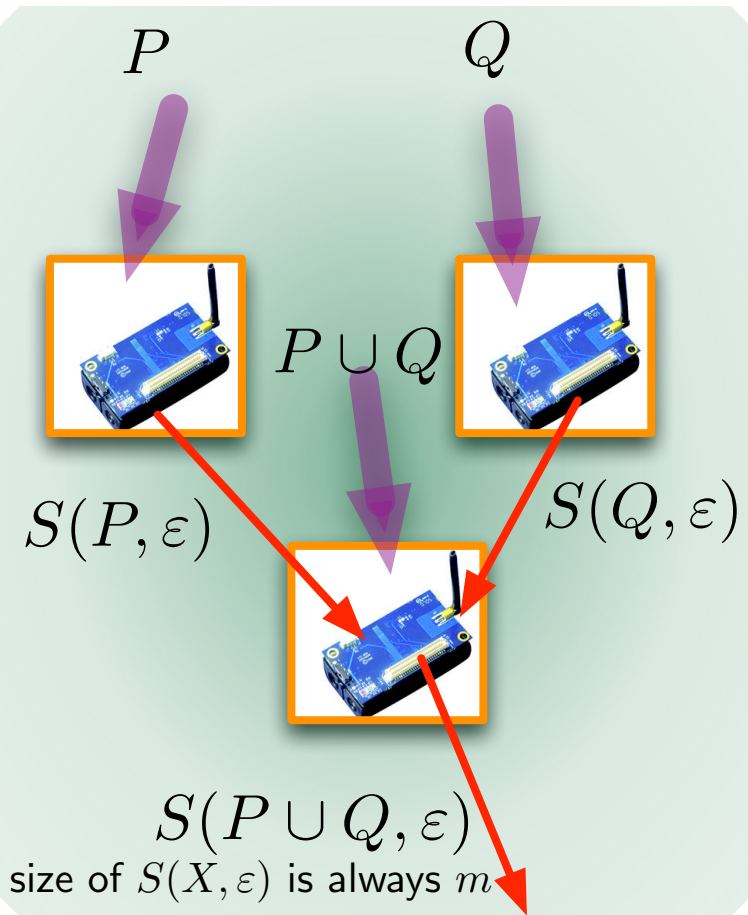
Massive Distributed Computation

data centers
sensor networks
multi-core



- similar to: MUD, Dremel
more restrictive, "natural"
- generalizes streaming
- archiving summaries

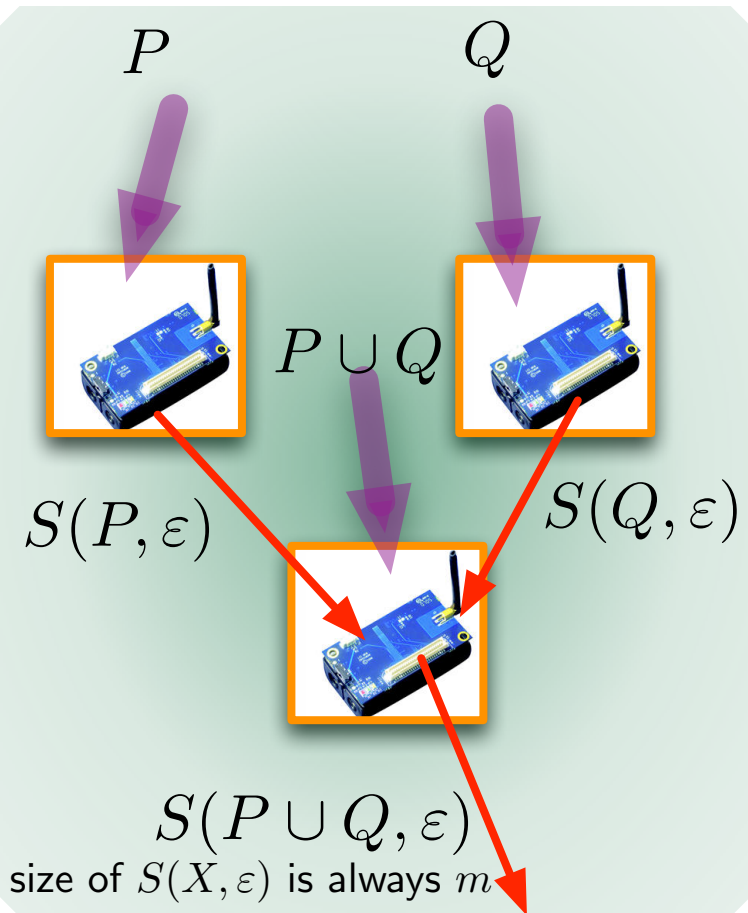
Random Sample



P

val	15	17	20	1	8	42	7	10	14	3
-----	----	----	----	---	---	----	---	----	----	---

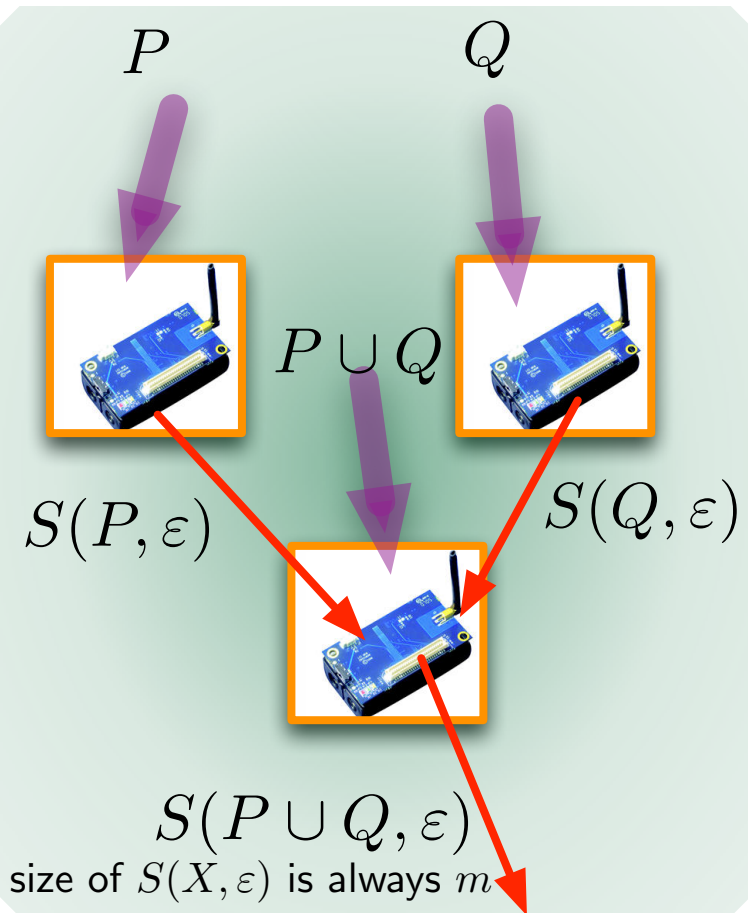
Random Sample



P

val	15	17	20	1	8	42	7	10	14	3
ran	.99	.42	.53	.01	.02	.23	.82	.75	.61	.14

Random Sample



P

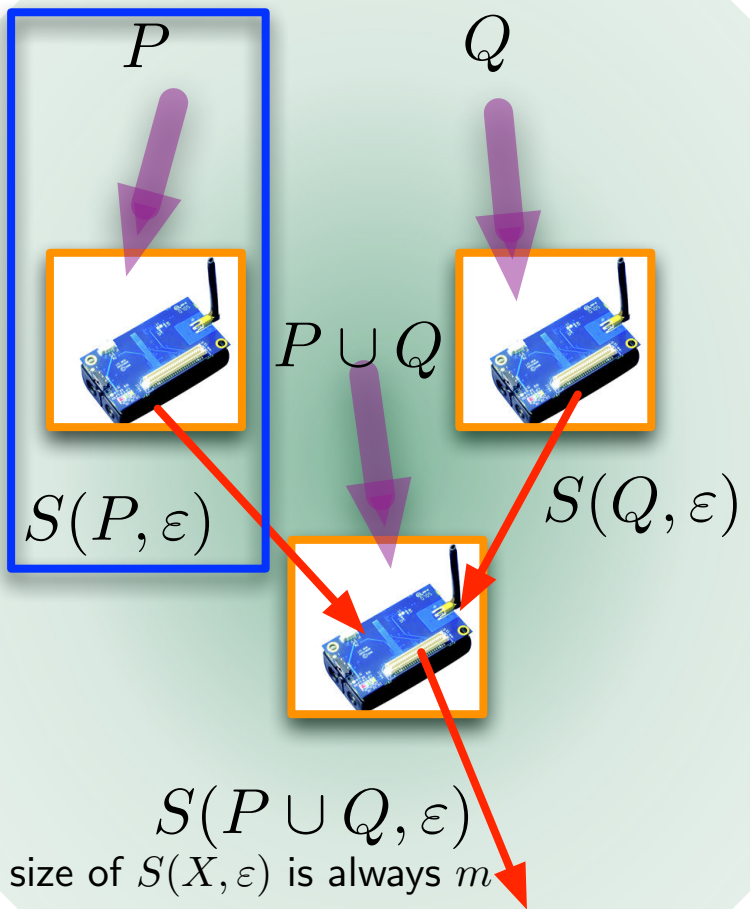
val	15	7	10	14	20	17	42	3	8	1
ran	.99	.82	.75	.61	.53	.42	.23	.14	.02	.01

Random Sample

$$S(P, \varepsilon)$$

P

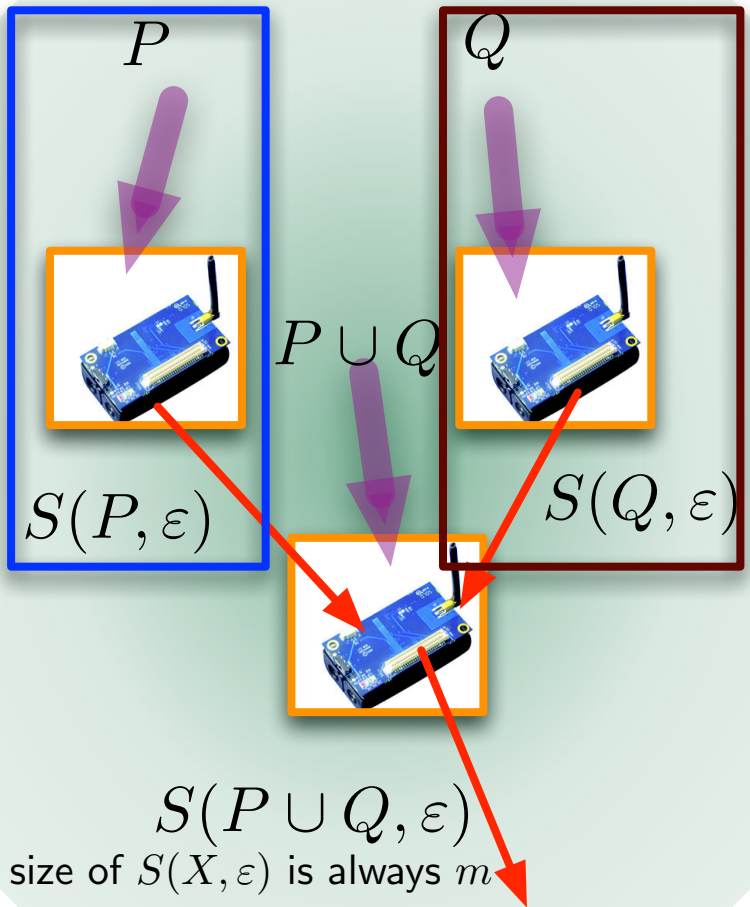
val	15	7	10	14	20	17	42	3	8	1
ran	.99	.82	.75	.61	.53	.42	.23	.14	.02	.01



$$S(P \cup Q, \varepsilon)$$

size of $S(X, \varepsilon)$ is always m

Random Sample



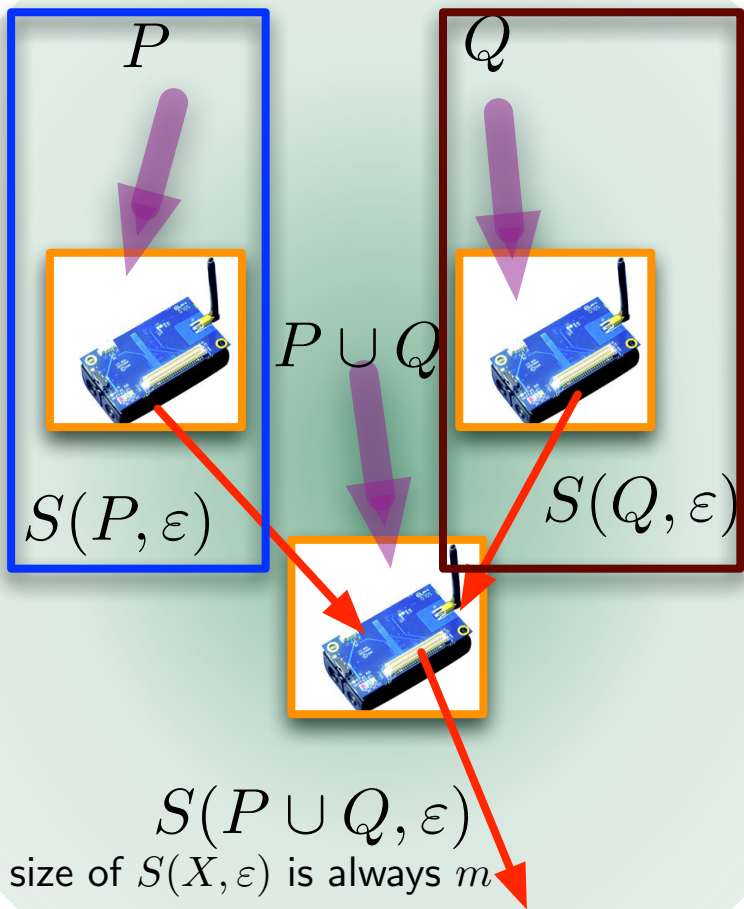
$S(P, \epsilon)$

P	val	15	7	10	14	20	17	42	3	8	1
	ran	.99	.82	.75	.61	.53	.42	.23	.14	.02	.01

$S(Q, \epsilon)$

Q	val	31	9	16	11	14	7	2	13	21	4
	ran	.90	.85	.80	.57	.50	.37	.31	.12	.10	.08

Random Sample



$S(P, \varepsilon)$

P

val	15	7	10	14	20	17	42	3	8	1
ran	.99	.82	.75	.61	.53	.42	.23	.14	.02	.01

$S(Q, \varepsilon)$

Q

val	31	9	16	11	14	7	2	13	21	4
ran	.90	.85	.80	.57	.50	.37	.31	.12	.10	.08

val	15	31	9	7	16	10
ran	.99	.90	.85	.82	.80	.75

Random Sample

$$S(P, \varepsilon)$$

P	val	15	7	10	14	20	17	42	3	8	1
	ran	.99	.82	.75	.61	.53	.42	.23	.14	.02	.01

$$S(Q, \varepsilon)$$

Q	val	31	9	16	11	14	7	2	13	21	4
	ran	.90	.85	.80	.57	.50	.37	.31	.12	.10	.08

$$S(P \cup Q, \varepsilon)$$

	val	15	31	9	7	16	10
	ran	.99	.90	.85	.82	.80	.75

$P \cup Q$



$$S(P \cup Q, \varepsilon)$$

size of $S(X, \varepsilon)$ is always m

Random Sample

$$S(P, \varepsilon)$$

P	val	15	7	10	14	20	17	42	3	8	1
	ran	.99	.82	.75	.61	.53	.42	.23	.14	.02	.01

$$S(Q, \varepsilon)$$

Q	val	31	9	16	11	14	7	2	13	21	4
	ran	.90	.85	.80	.57	.50	.37	.31	.12	.10	.08

$$S(P \cup Q, \varepsilon)$$

$P \cup Q$	val	15	31	9	7	16	10
	ran	.99	.90	.85	.82	.80	.75

$P \cup Q$



$$S(P \cup Q, \varepsilon)$$

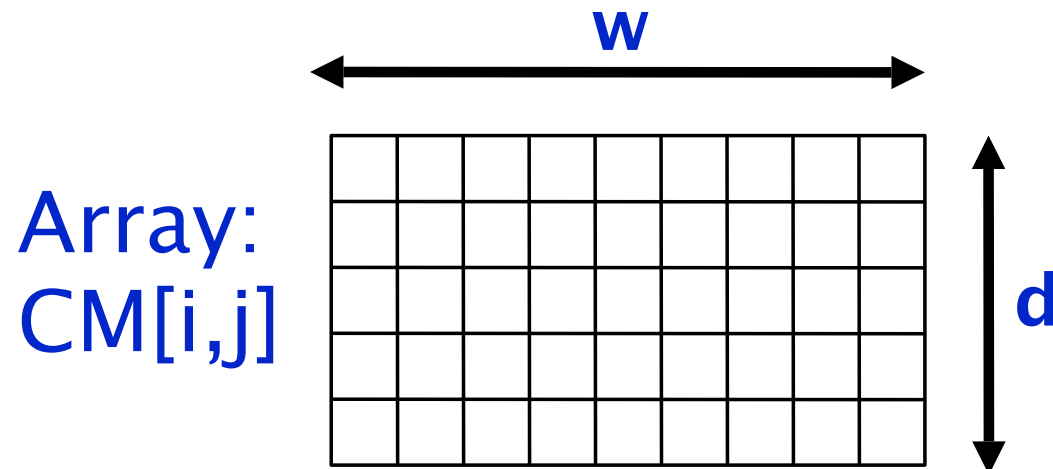
size of $S(X, \varepsilon)$ is always m

max element
top k elements

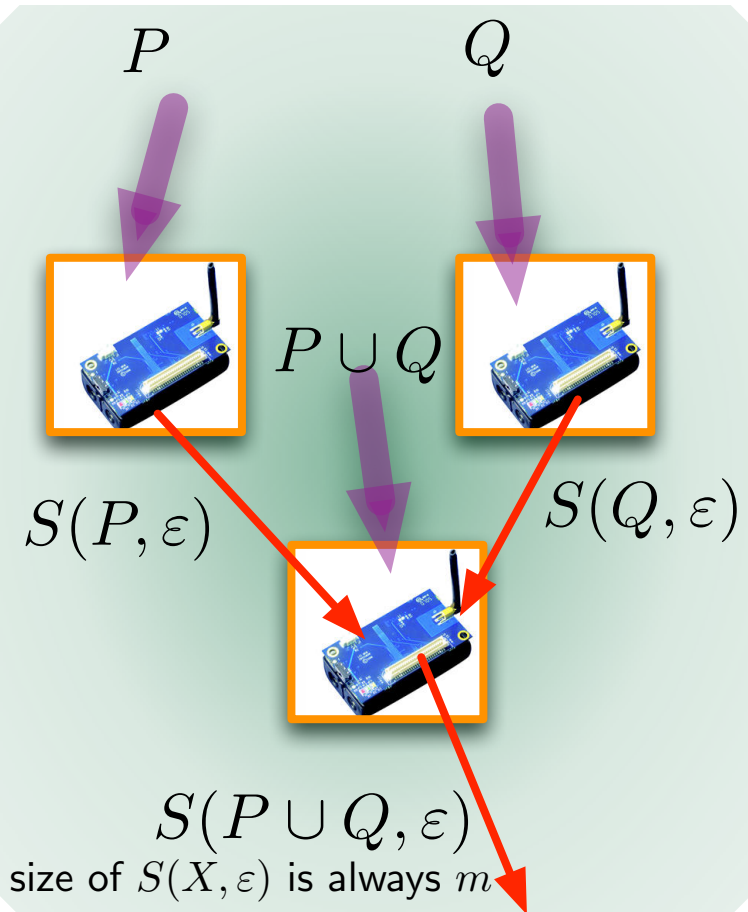
Linear Sketches

Count-Min sketch of vector $P[1..U]$:

- Linear sketch as array size $w \times d$
- Use d hash functions h to map x to $[1..w]$
- Estimate $P[i] = \min_j \text{CM}[h_j(i), j]$



Linear Sketches

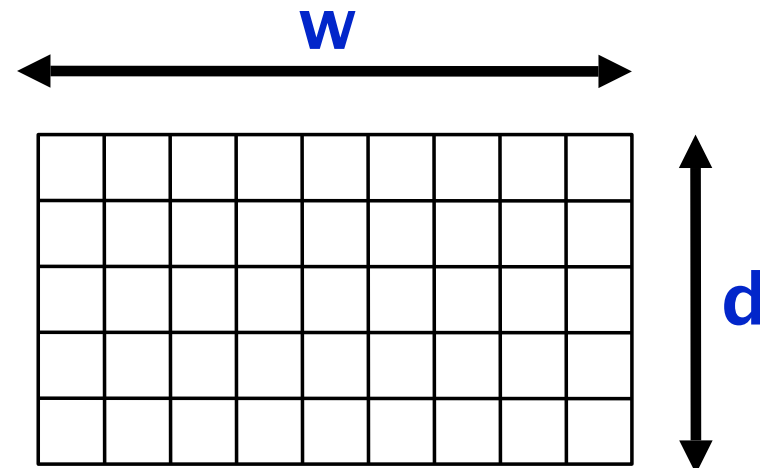


Count-Min sketch of vector $P[1..U]$:

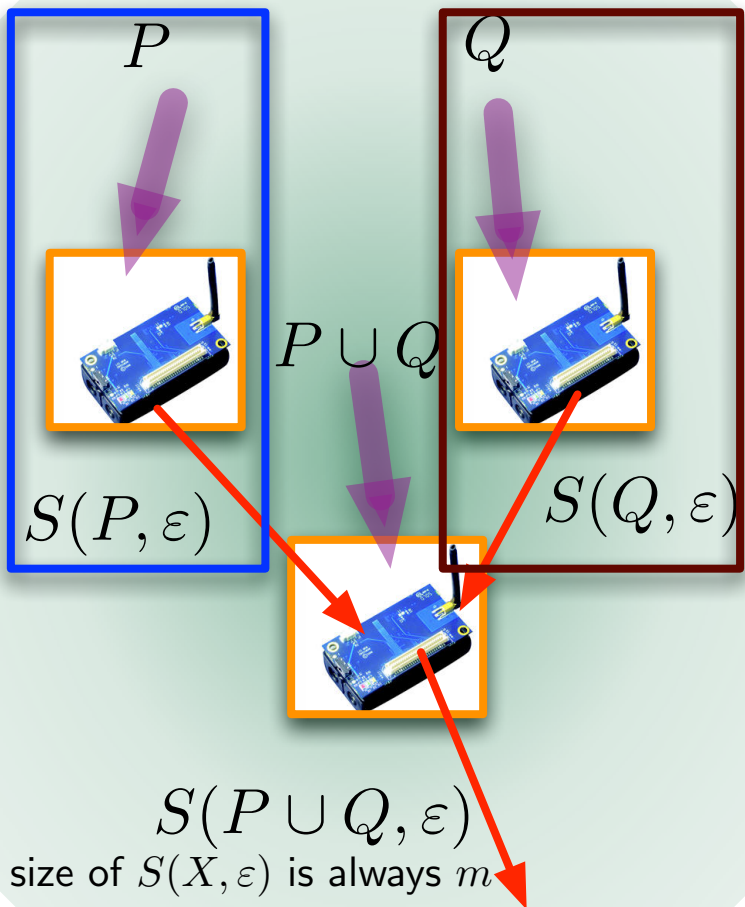
- Linear sketch as array size $w \times d$
- Use d hash functions h to map x to $[1..w]$
- Estimate $P[i] = \min_j \text{CM}[h_j(i), j]$

Mergeable: $\text{CM}(P + Q) = \text{CM}(P) + \text{CM}(Q)$

Array:
 $\text{CM}[i,j]$



Linear Sketches



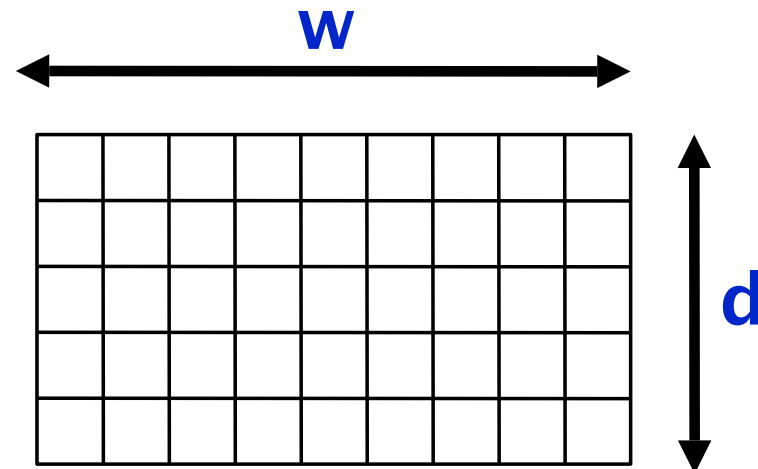
Count-Min sketch of vector $P[1..U]$:

- Linear sketch as array size $w \times d$
- Use d hash functions h to map x to $[1..w]$
- Estimate $P[i] = \min_j \text{CM}[h_j(i), j]$

Mergeable: $\text{CM}(P + Q) = \text{CM}(P) + \text{CM}(Q)$

$S(P, \epsilon)$ $S(Q, \epsilon)$

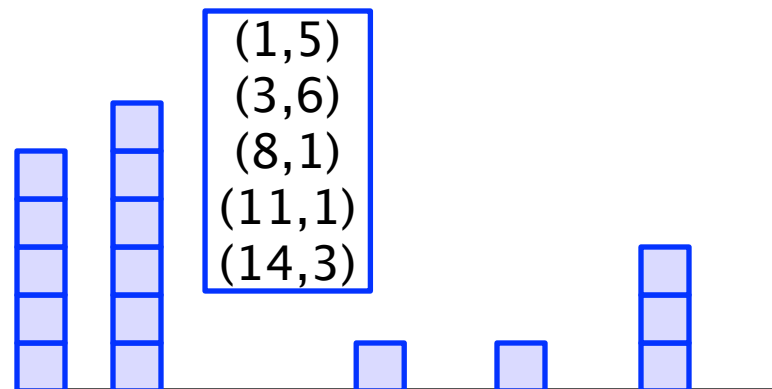
Array:
 $\text{CM}[i,j]$



Heavy Hitters Summaries

Misra-Gries (MG) sketch of $P[1..U]$:

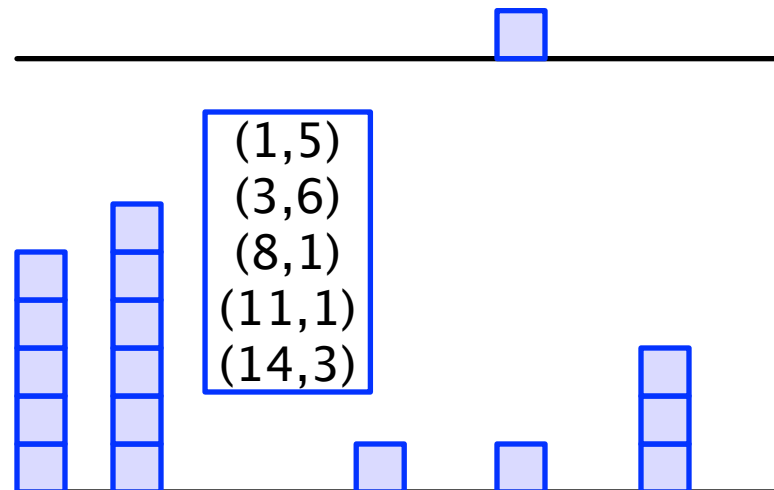
- Keep k (index,count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair, or decrement all counts



Heavy Hitters Summaries

Misra-Gries (MG) sketch of $P[1..U]$:

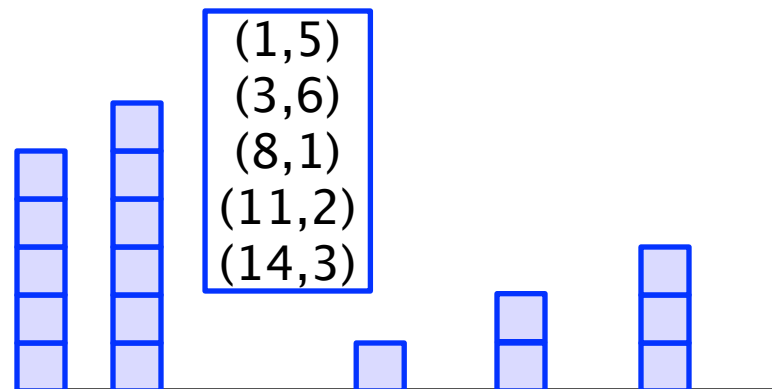
- Keep k (index,count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair,
or decrement all counts



Heavy Hitters Summaries

Misra-Gries (MG) sketch of $P[1..U]$:

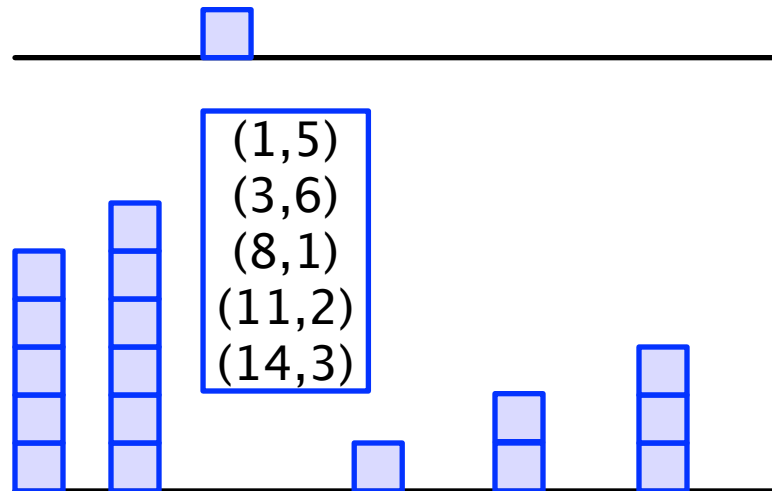
- Keep k (index,count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair,
or decrement all counts



Heavy Hitters Summaries

Misra-Gries (MG) sketch of $P[1..U]$:

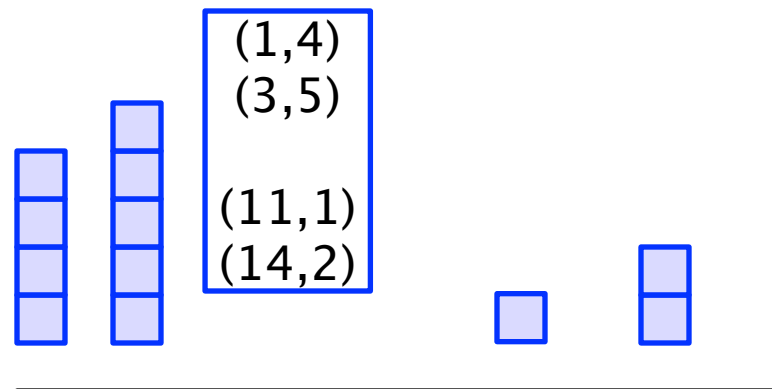
- Keep k (index,count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair,
or decrement all counts



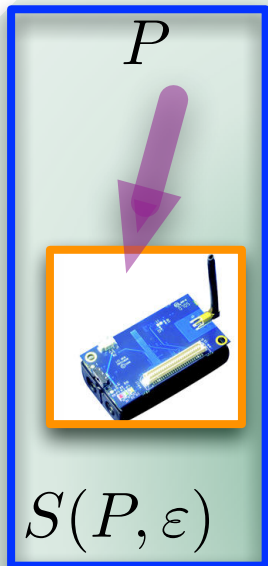
Heavy Hitters Summaries

Misra-Gries (MG) sketch of $P[1..U]$:

- Keep k (index,count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair,
or decrement all counts



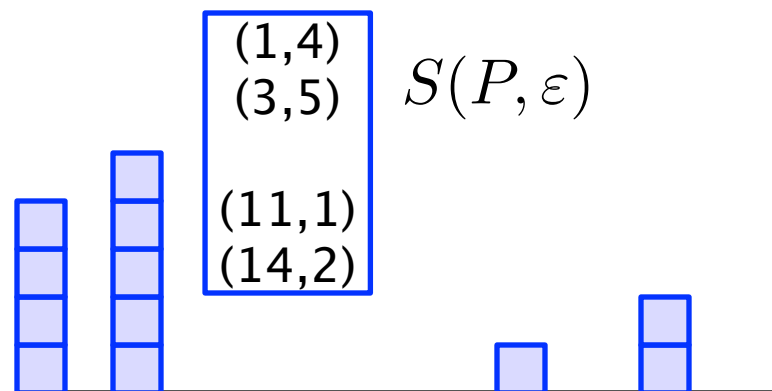
Heavy Hitters Summaries



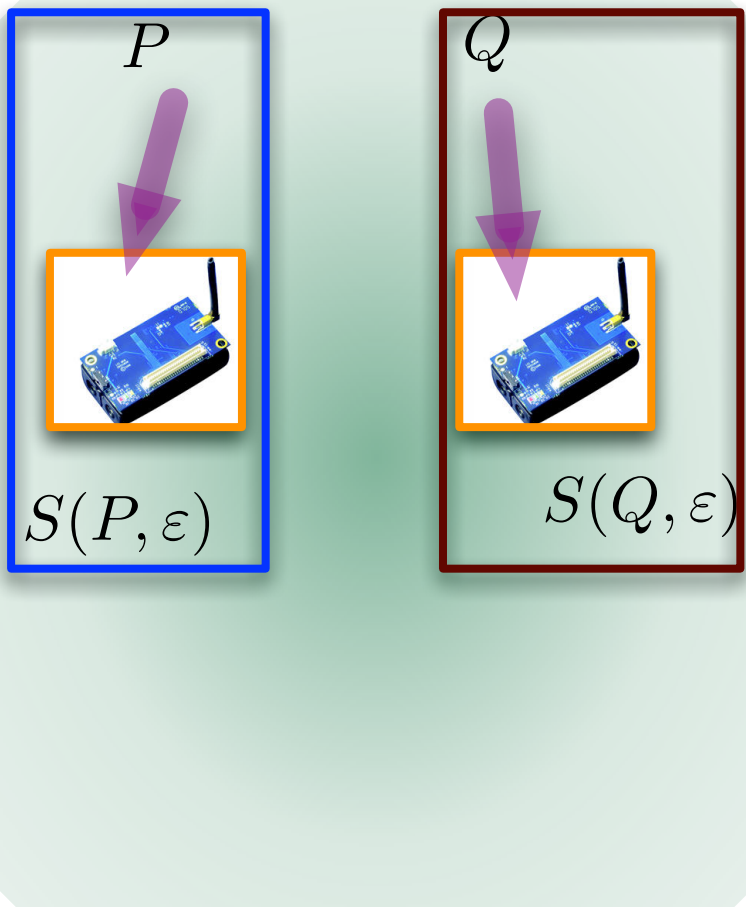
Misra-Gries (MG) sketch of $P[1..U]$:

- Keep k (index, count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair, or decrement all counts

$$|P[i] - \text{MG}[i]| \leq \epsilon = \hat{m}/(k + 1)$$



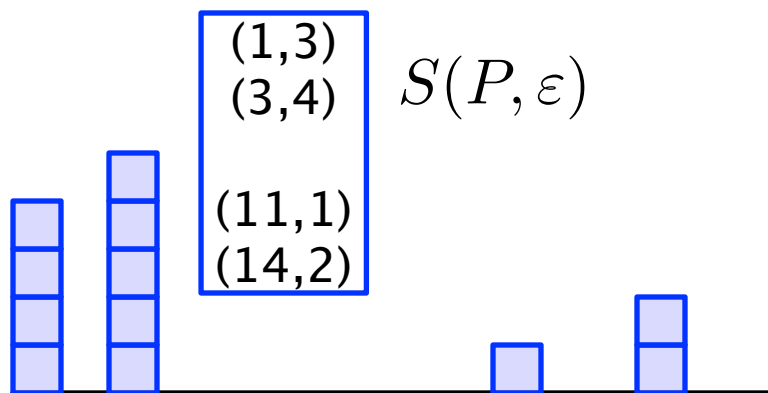
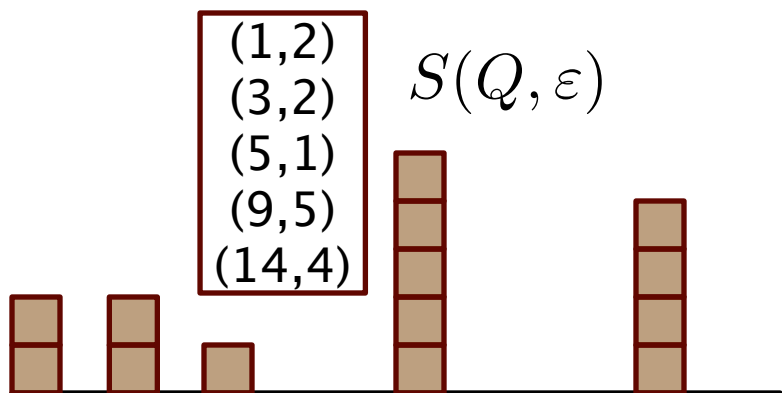
Heavy Hitters Summaries



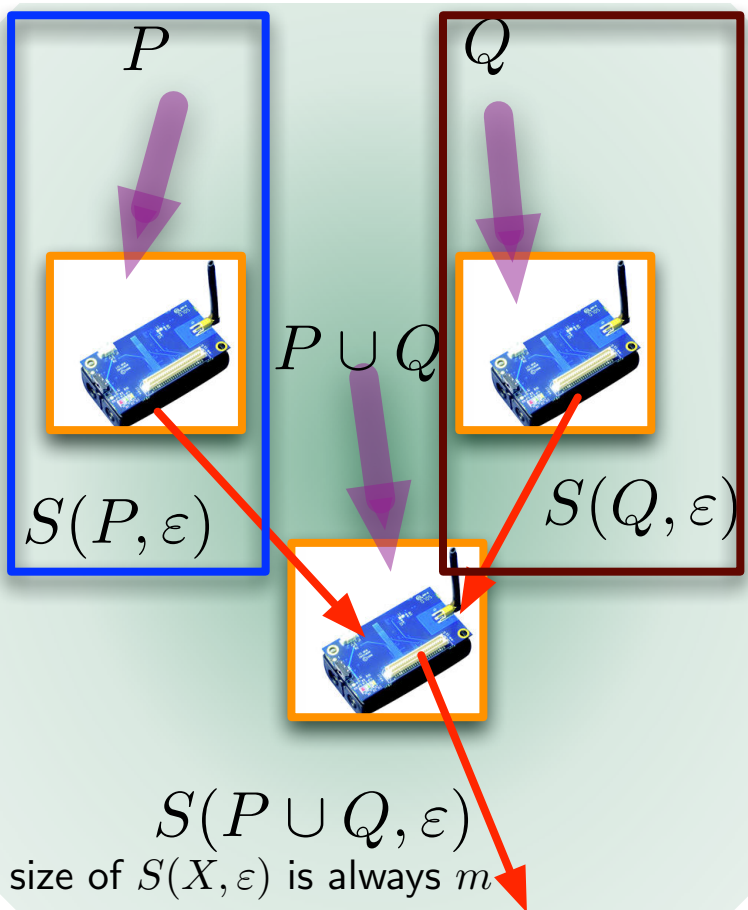
Misra-Gries (MG) sketch of $P[1..U]$:

- Keep k (index, count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair, or decrement all counts

Mergeable: Stack $MG(P) + MG(Q)$,



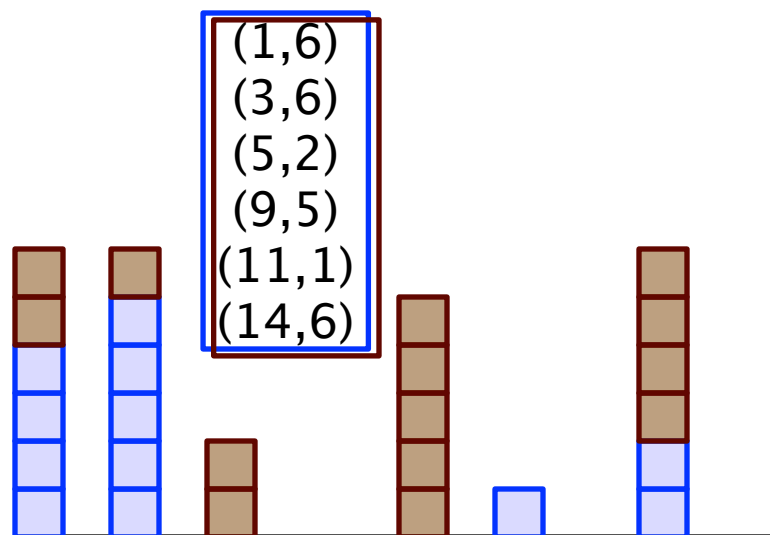
Heavy Hitters Summaries



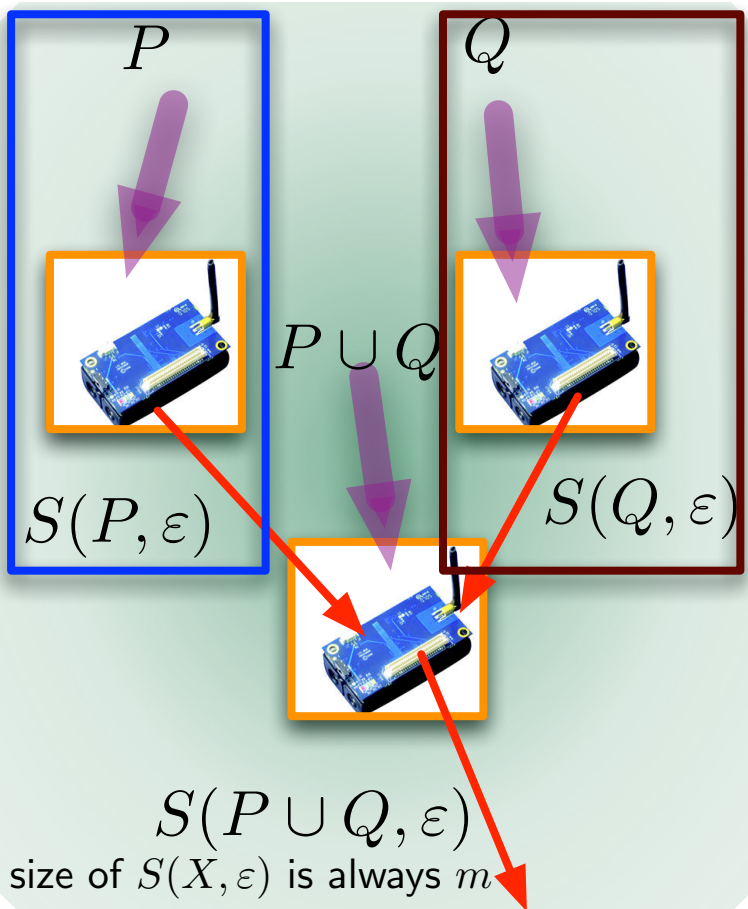
Misra-Gries (MG) sketch of $P[1..U]$:

- Keep k (index, count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair, or decrement all counts

Mergeable: Stack $MG(P) + MG(Q)$,



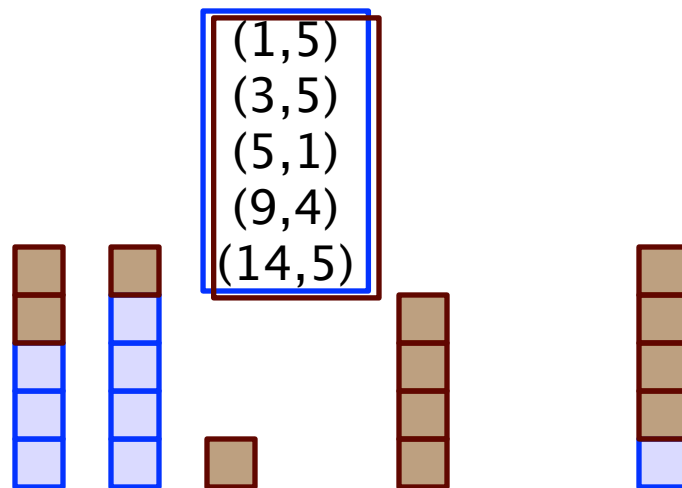
Heavy Hitters Summaries



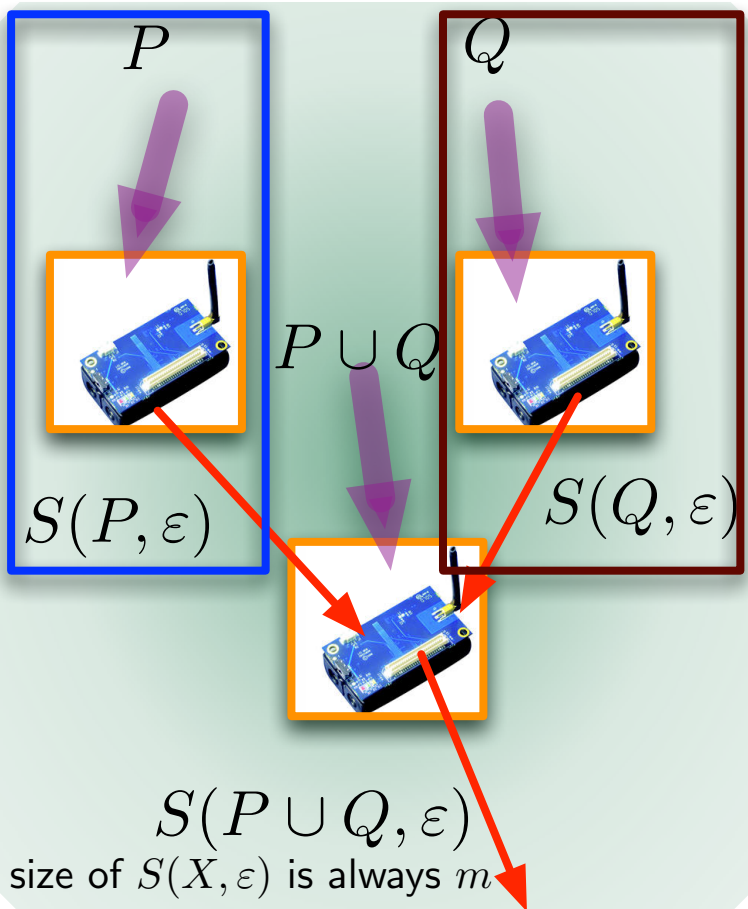
Misra-Gries (MG) sketch of $P[1..U]$:

- Keep k (index, count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair, or decrement all counts

Mergeable: Stack $MG(P) + MG(Q)$,
decrement all counts C_{k+1}



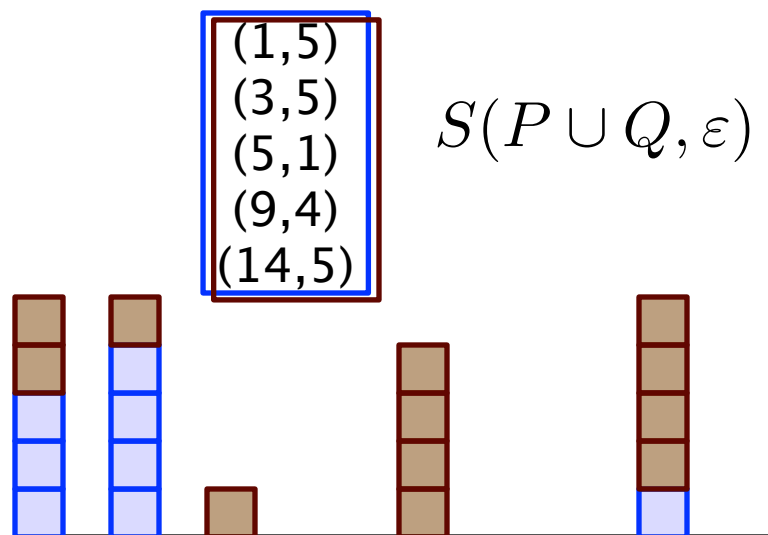
Heavy Hitters Summaries



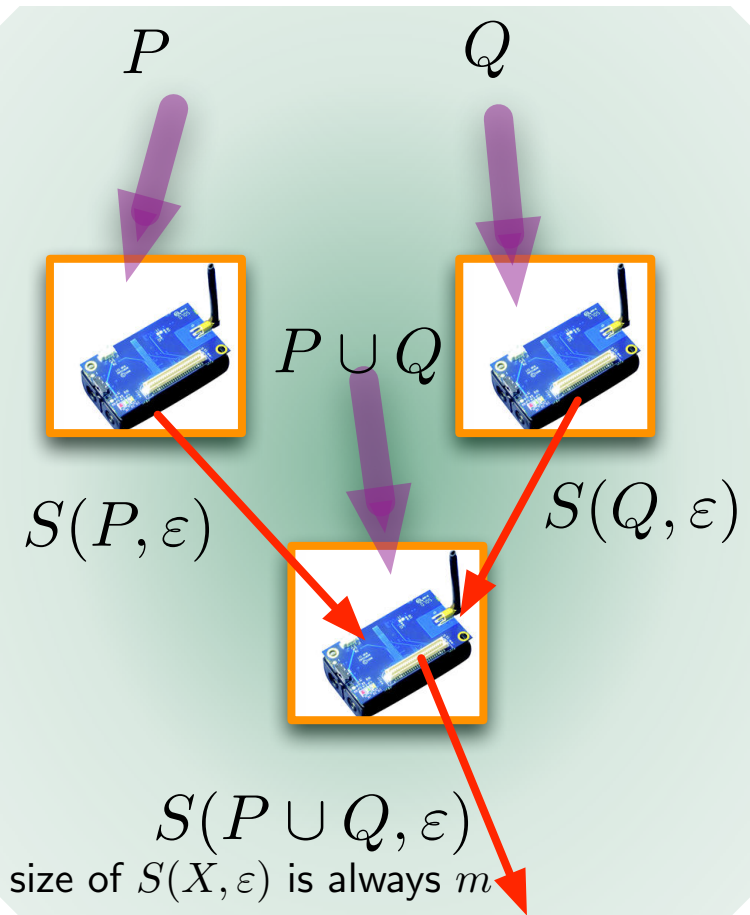
Misra-Gries (MG) sketch of $P[1..U]$:

- Keep k (index, count) pairs
- If existing index arrives, update count
- If new index arrives, make new pair, or decrement all counts

Mergeable: Stack $MG(P) + MG(Q)$,
decrement all counts C_{k+1}

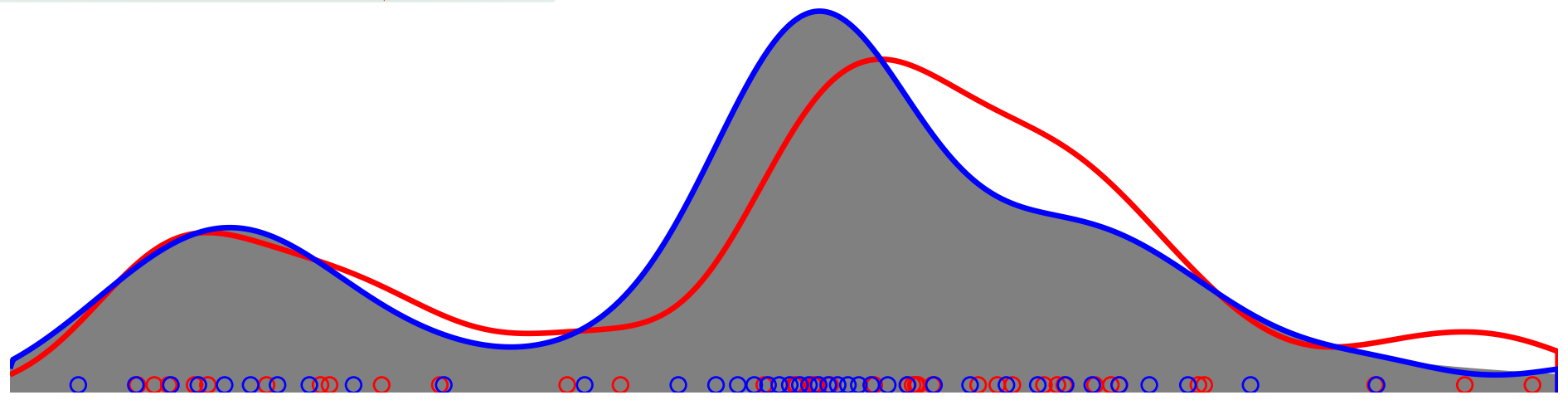


ϵ -Samples (Intervals)

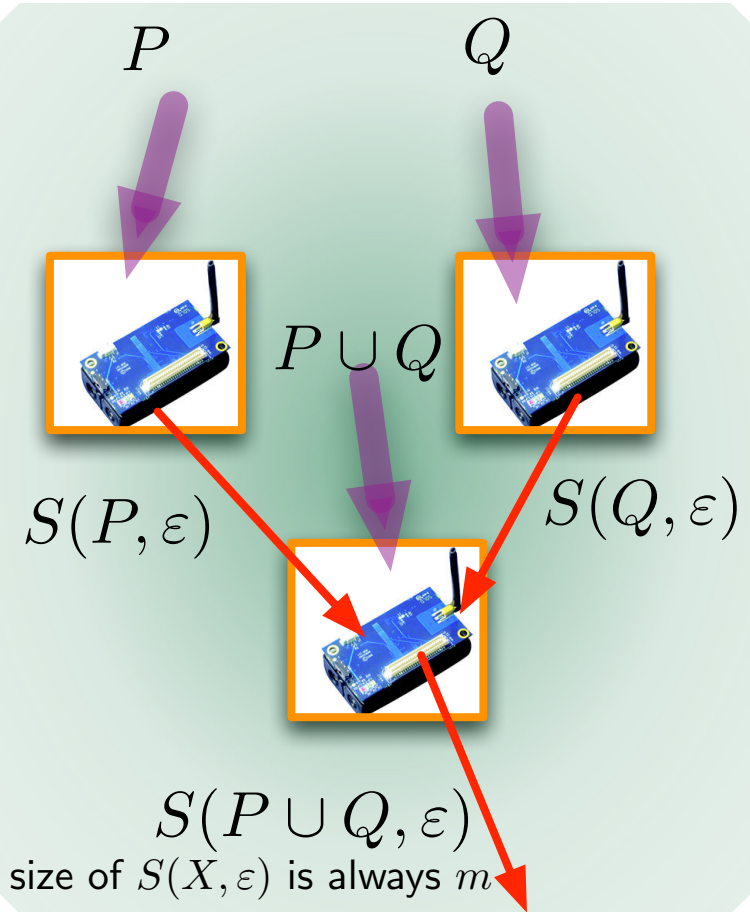


P

val	15	17	20	1	8	42	7	10	14	3
-----	----	----	----	---	---	----	---	----	----	---

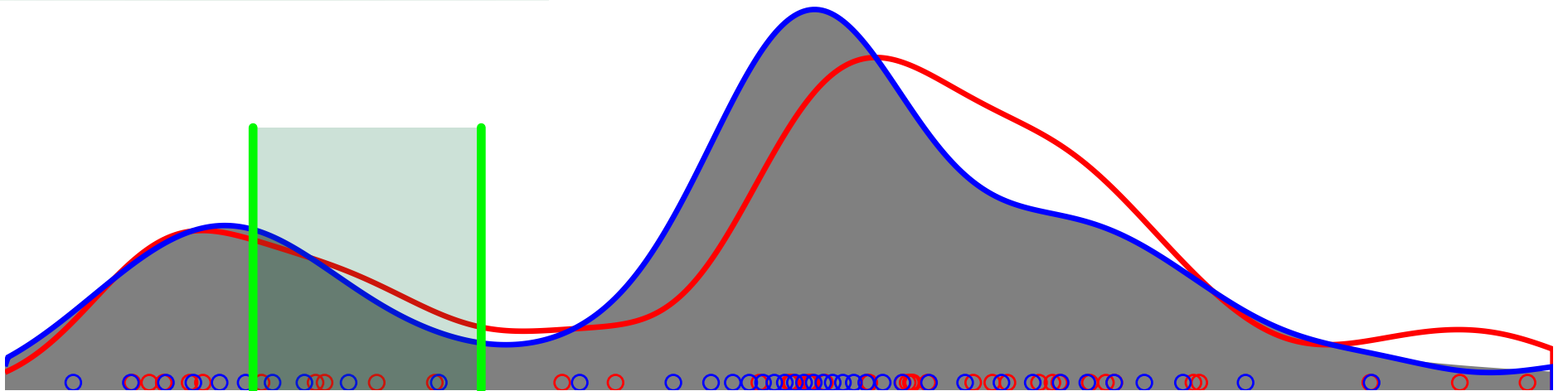


ϵ -Samples (Intervals)

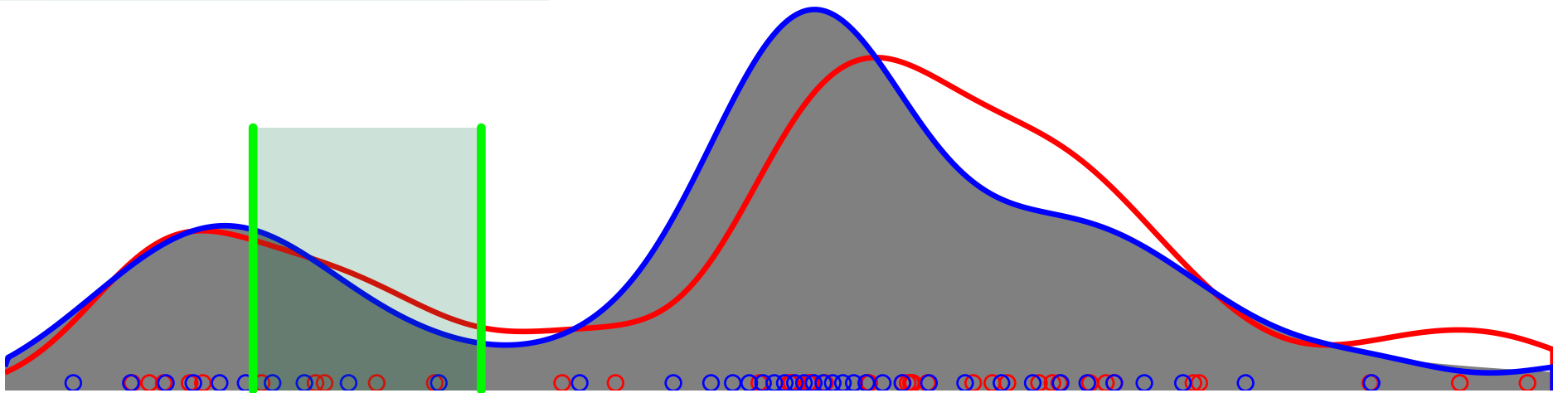
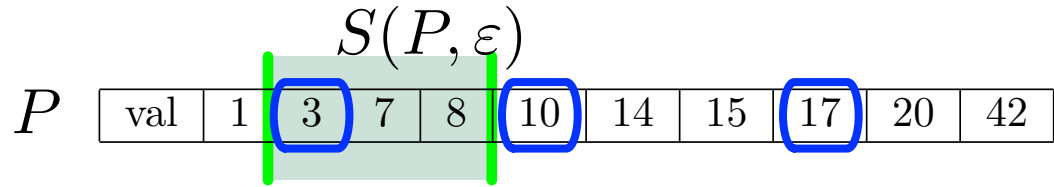
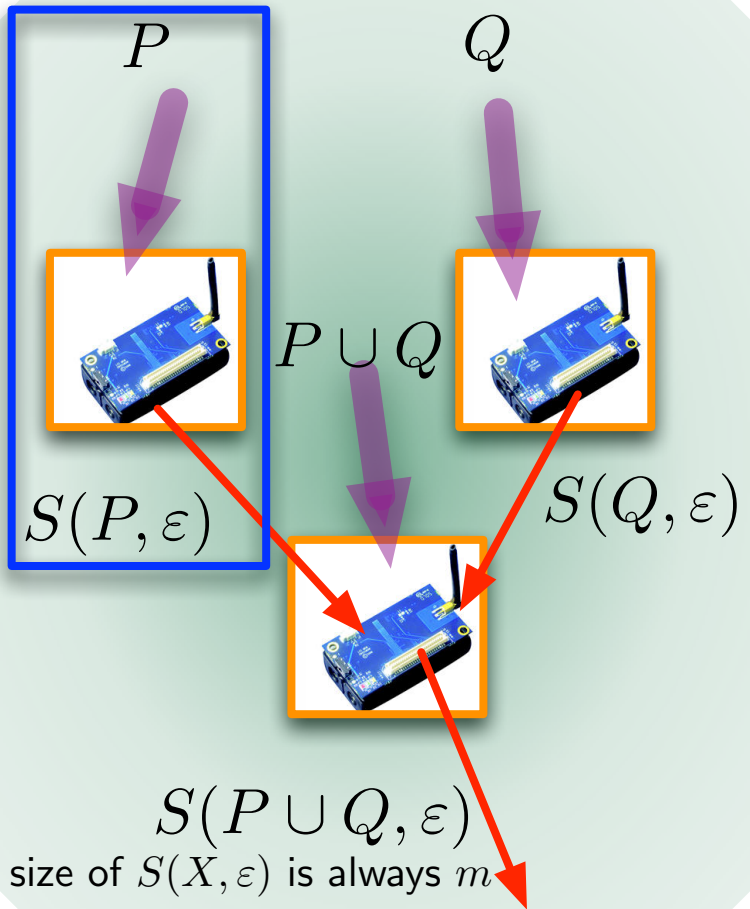


P

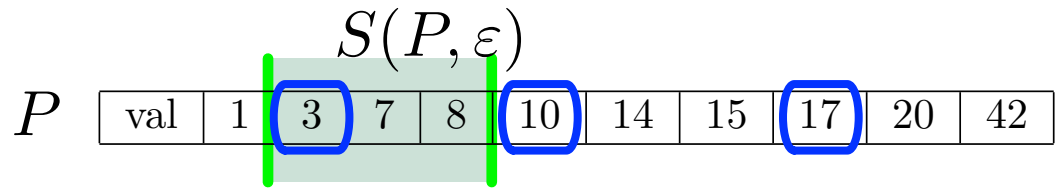
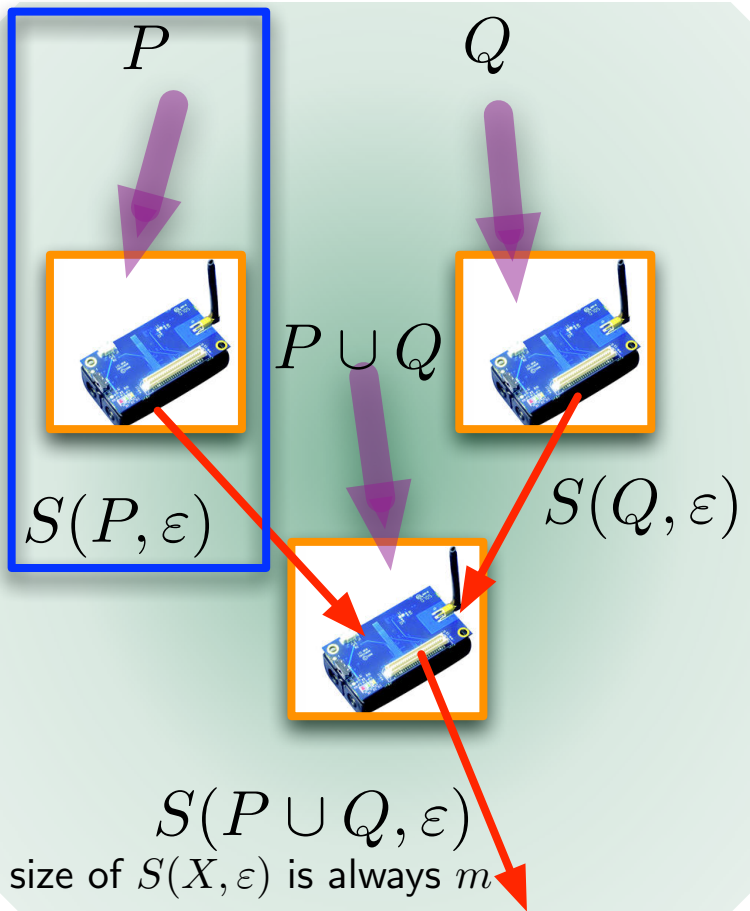
val	15	17	20	1	8	42	7	10	14	3
-----	----	----	----	---	---	----	---	----	----	---



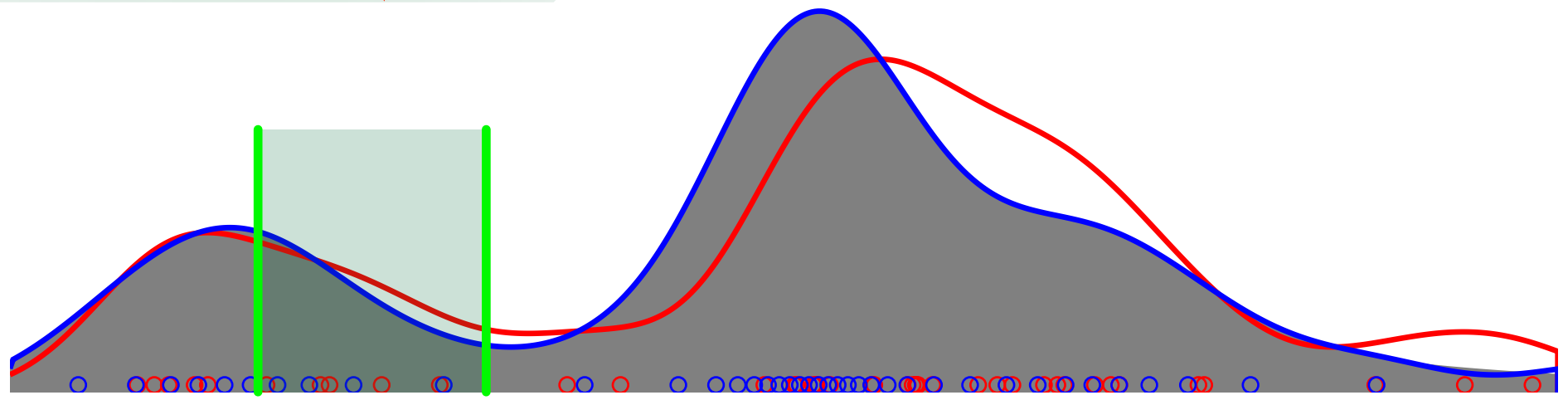
ϵ -Samples (Intervals)



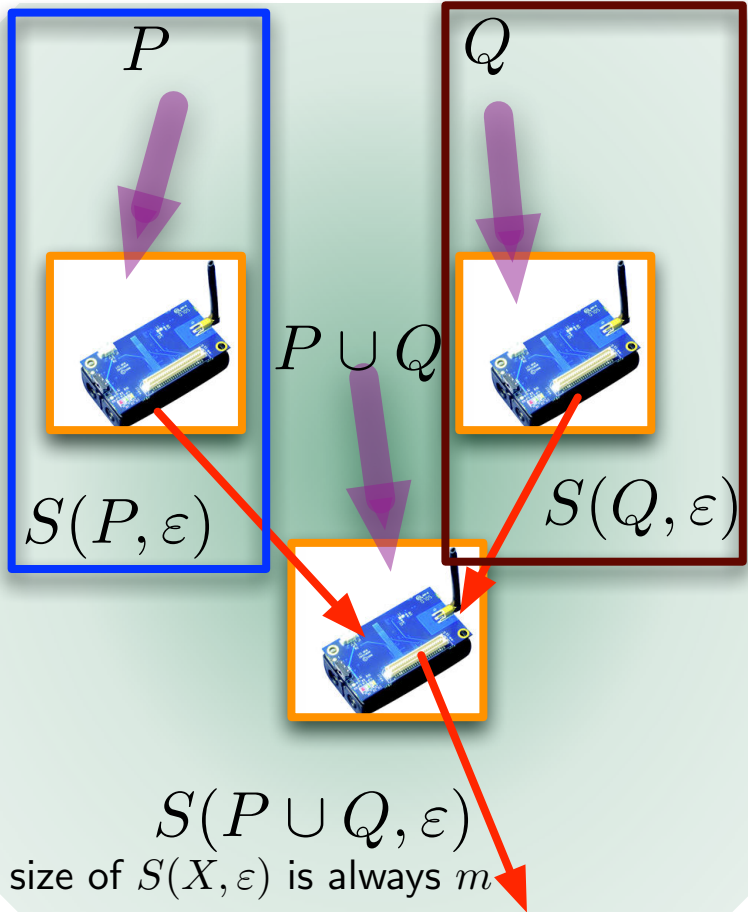
ϵ -Samples (Intervals)



An ϵ -sample of ϵ -sample is a 2ϵ -sample



ϵ -Samples (Intervals)



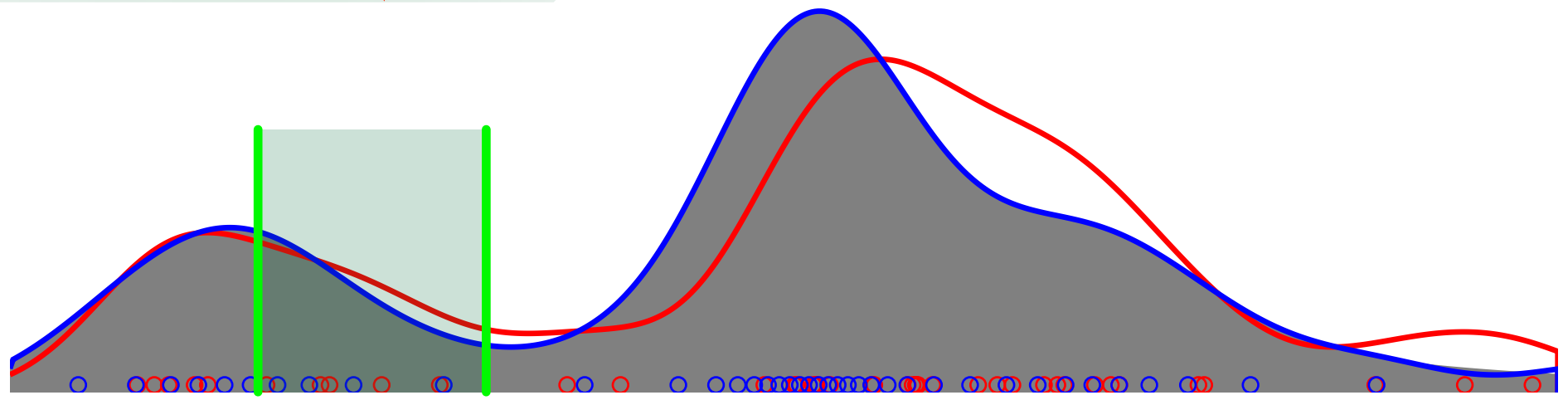
$S(P, \epsilon)$

P	val	1	3	7	8	10	14	15	17	20	42
-----	-----	---	---	---	---	----	----	----	----	----	----

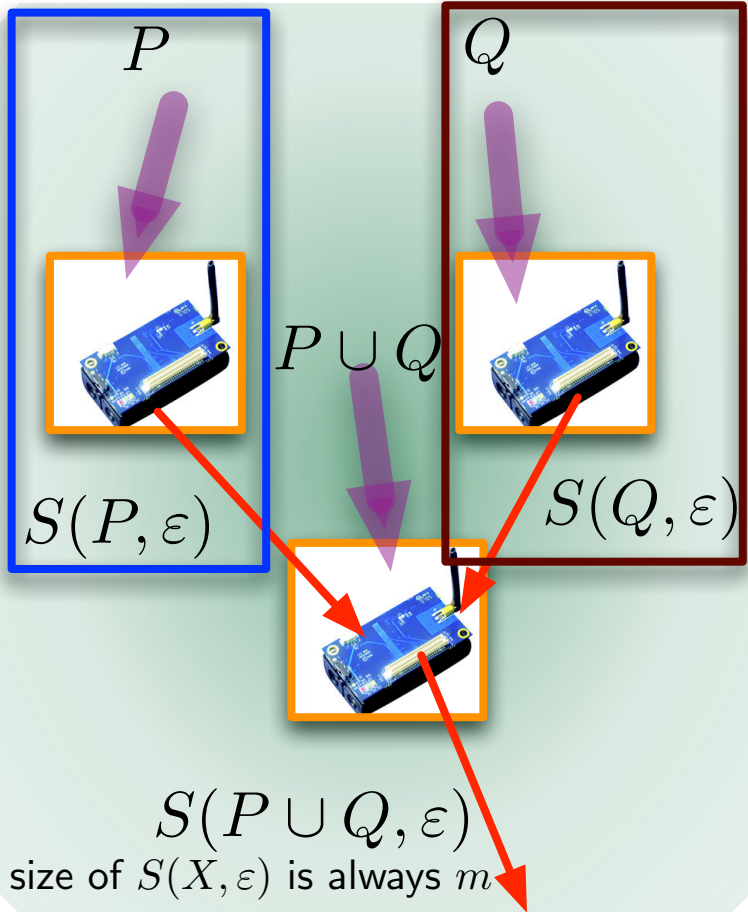
$S(Q, \epsilon)$

Q	val	2	4	7	9	11	13	14	16	21	31
-----	-----	---	---	---	---	----	----	----	----	----	----

val	3	4	10	11	16	17
-----	---	---	----	----	----	----



ϵ -Samples (Intervals)



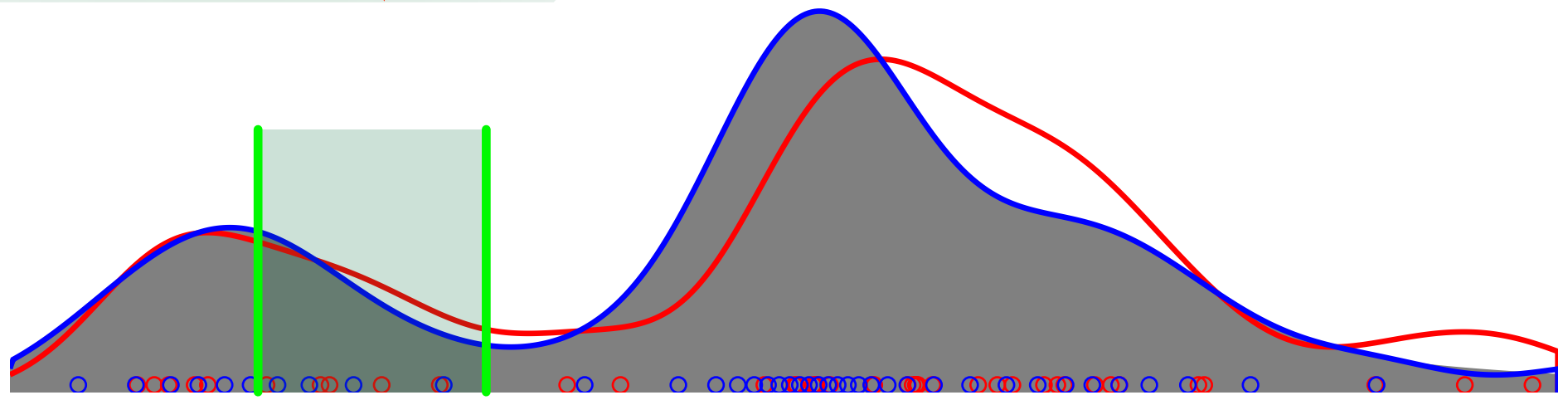
$S(P, \epsilon)$

P	val	1	3	7	8	10	14	15	17	20	42
-----	-----	---	---	---	---	----	----	----	----	----	----

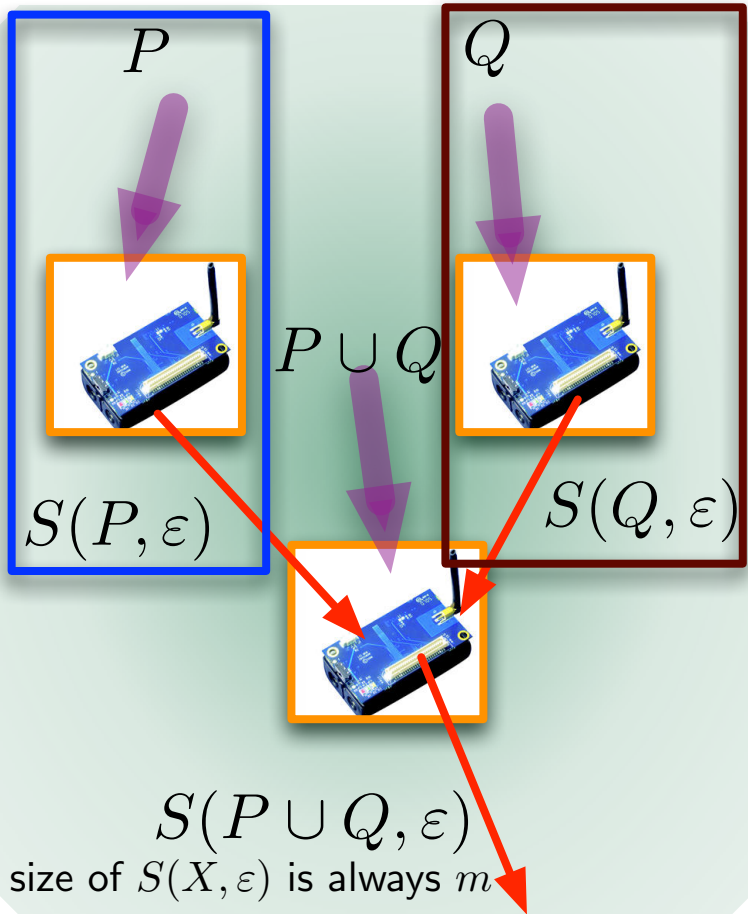
$S(Q, \epsilon)$

Q	val	2	4	7	9	11	13	14	16	21	31
-----	-----	---	---	---	---	----	----	----	----	----	----

val	3	4	10	11	16	17
-----	---	---	----	----	----	----



ϵ -Samples (Intervals)



$S(P, \epsilon)$

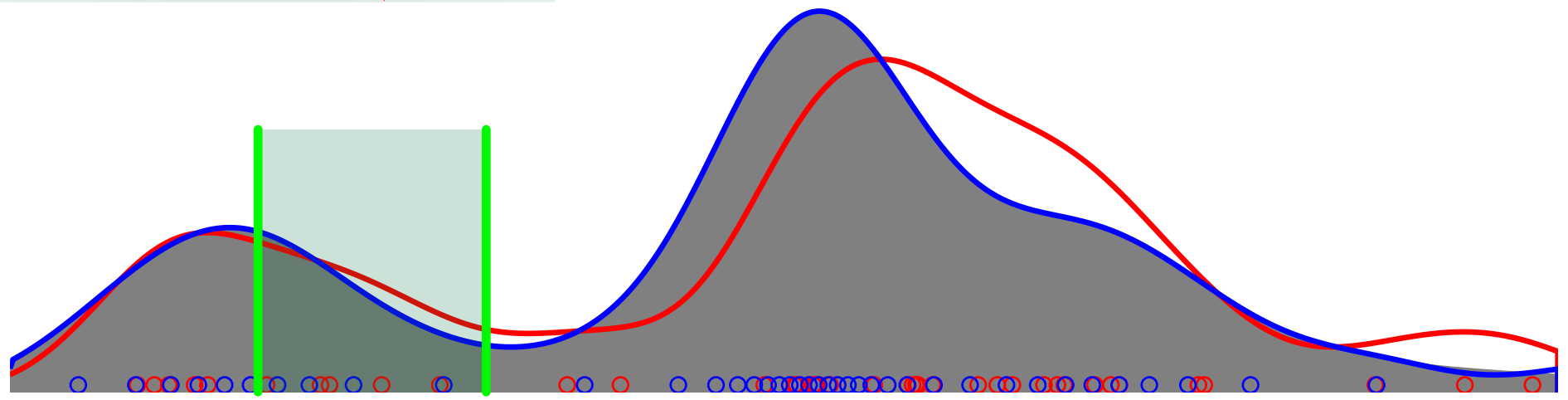
P	val	1	3	7	8	10	14	15	17	20	42
-----	-----	---	---	---	---	----	----	----	----	----	----

$S(Q, \epsilon)$

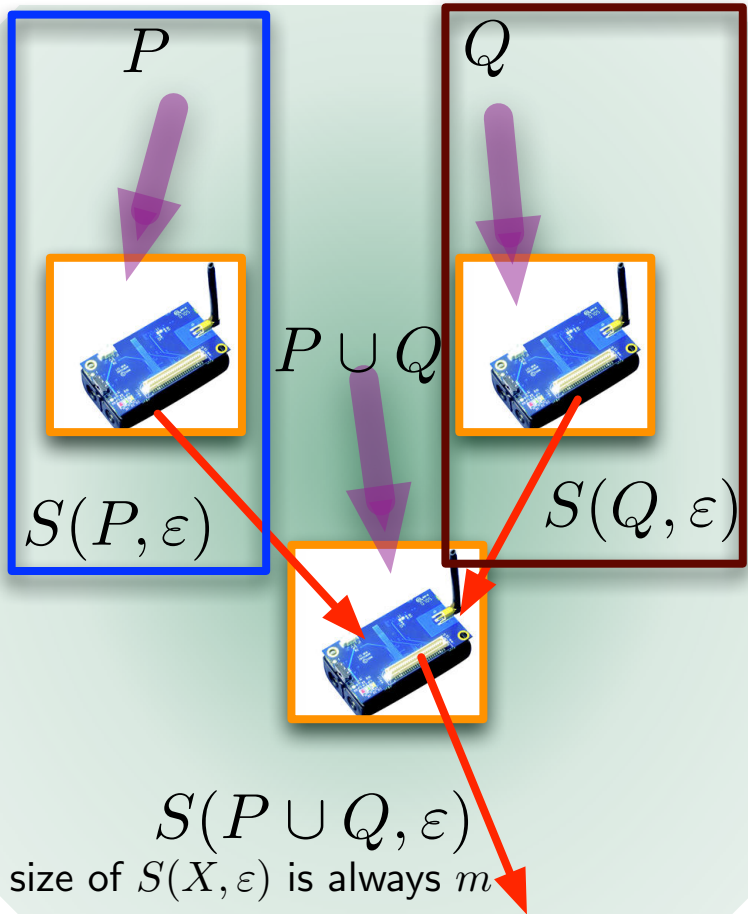
Q	val	2	4	7	9	11	13	14	16	21	31
-----	-----	---	---	---	---	----	----	----	----	----	----

$S(P \cup Q, \epsilon)$

	val	3	4	10	11	16	17
--	-----	---	---	----	----	----	----



ϵ -Samples (Intervals)



$S(P, \epsilon)$

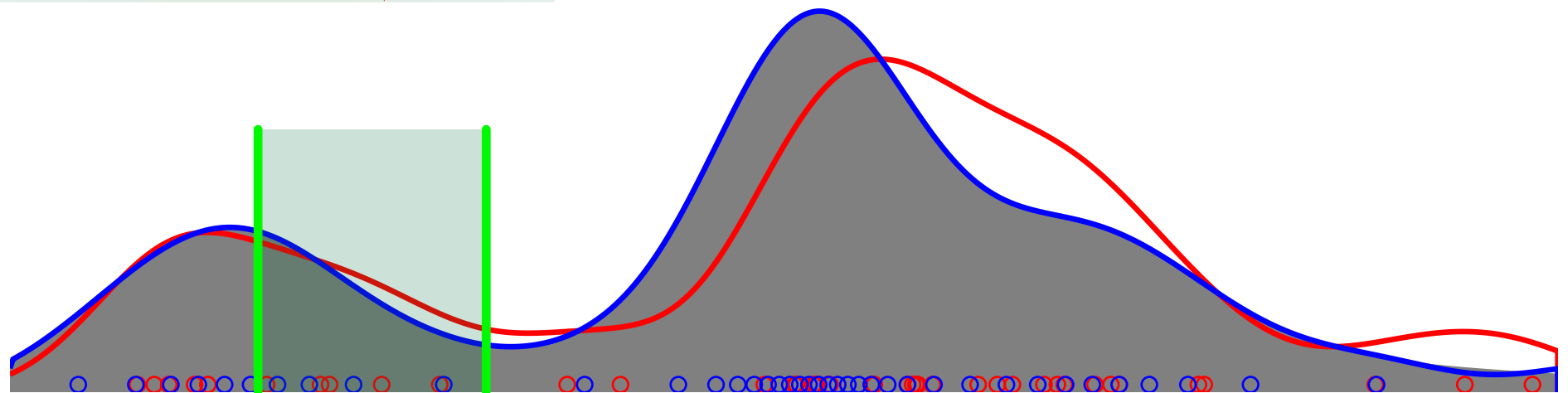
P	val	1	3	7	8	10	14	15	17	20	42
-----	-----	---	---	---	---	----	----	----	----	----	----

$S(Q, \epsilon)$

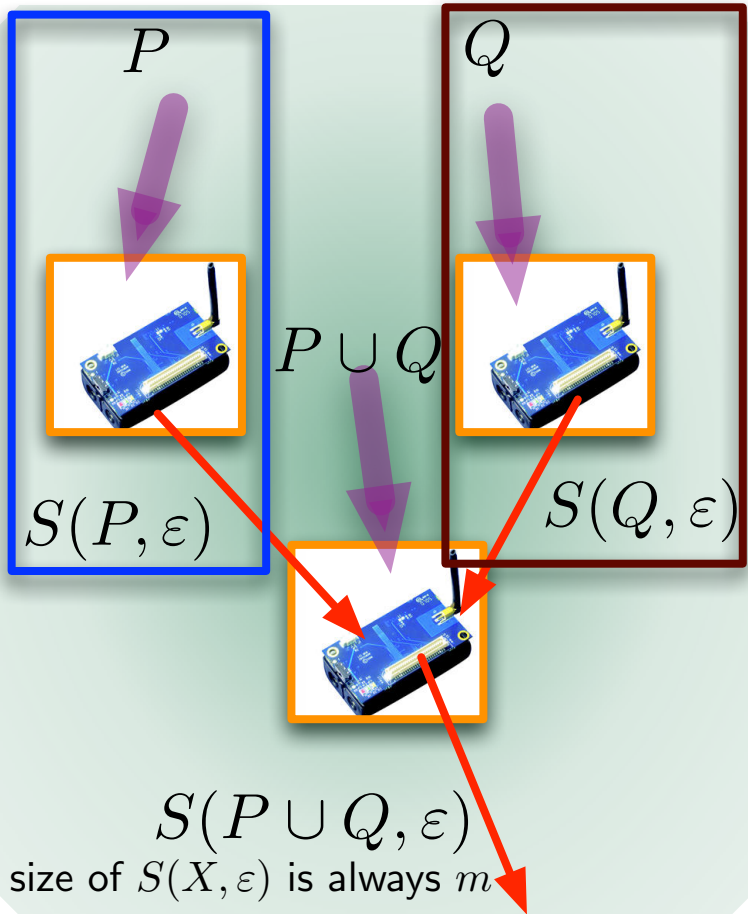
Q	val	2	4	7	9	11	13	14	16	21	31
-----	-----	---	---	---	---	----	----	----	----	----	----

$S(P \cup Q, \epsilon)$

	val	3	4	10	11	16	17
--	-----	---	---	----	----	----	----



ϵ -Samples (Intervals)



$S(P, \epsilon)$

P	val	1	3	7	8	10	14	15	17	20	42
-----	-----	---	---	---	---	----	----	----	----	----	----

$S(Q, \epsilon)$

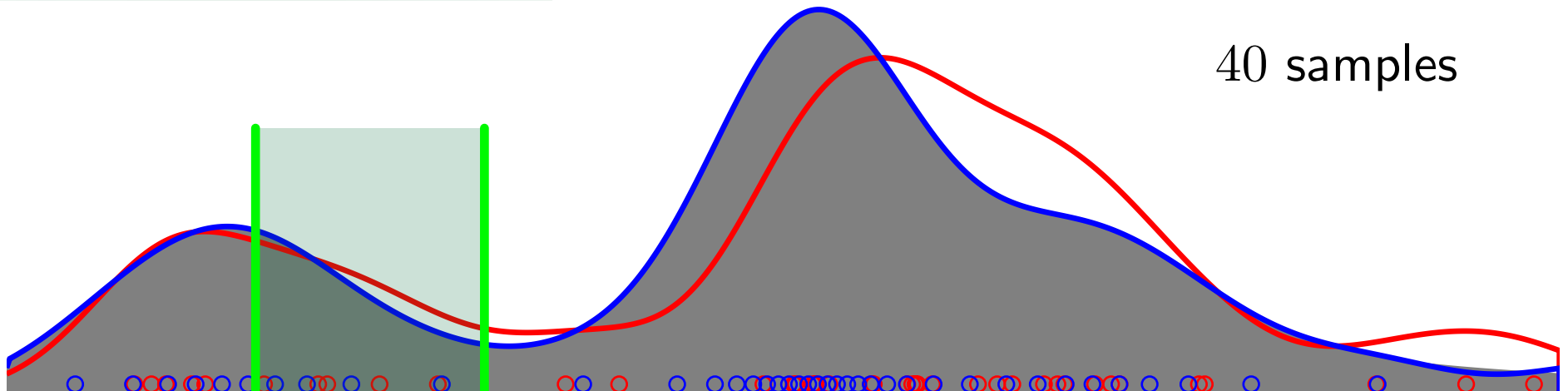
Q	val	2	4	7	9	11	13	14	16	21	31
-----	-----	---	---	---	---	----	----	----	----	----	----

$S(P \cup Q, \epsilon)$

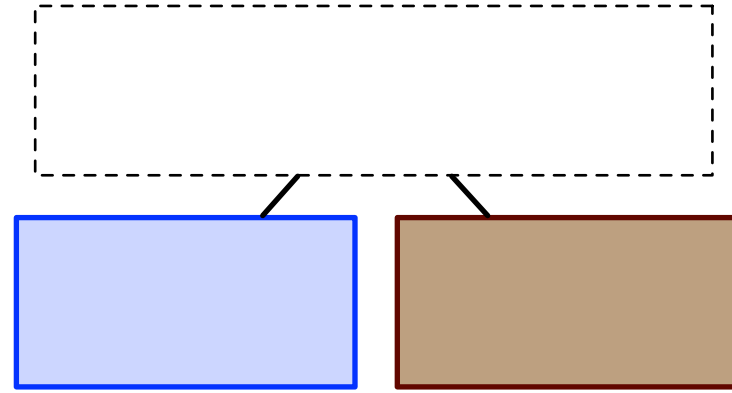
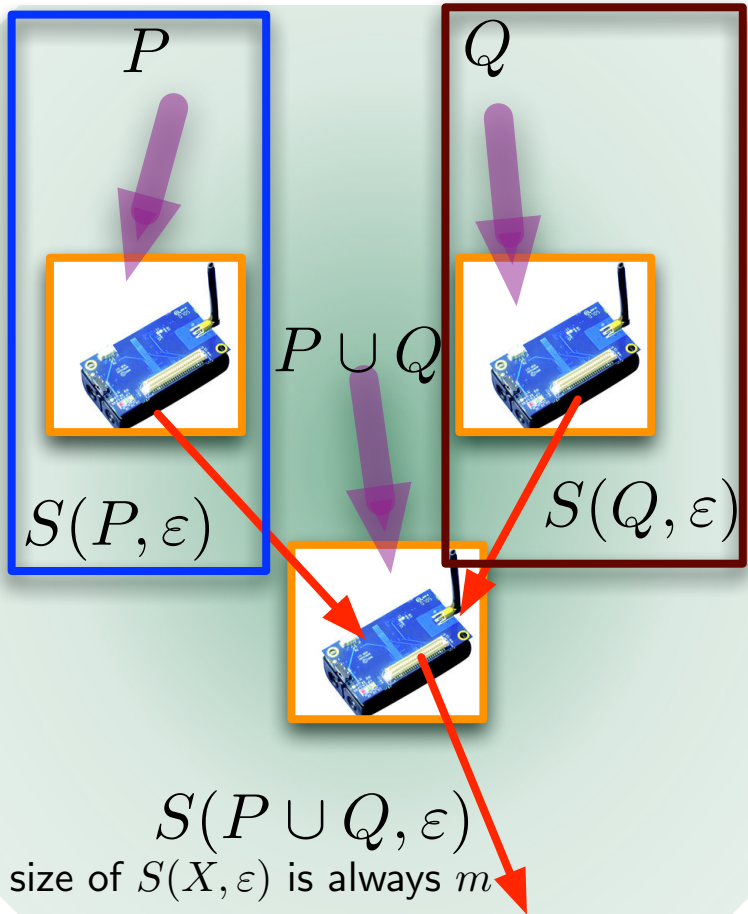
	val	3	4	10	11	16	17
--	-----	---	---	----	----	----	----

Random Sample: $(1/\epsilon^2) \log(1/\delta)$.

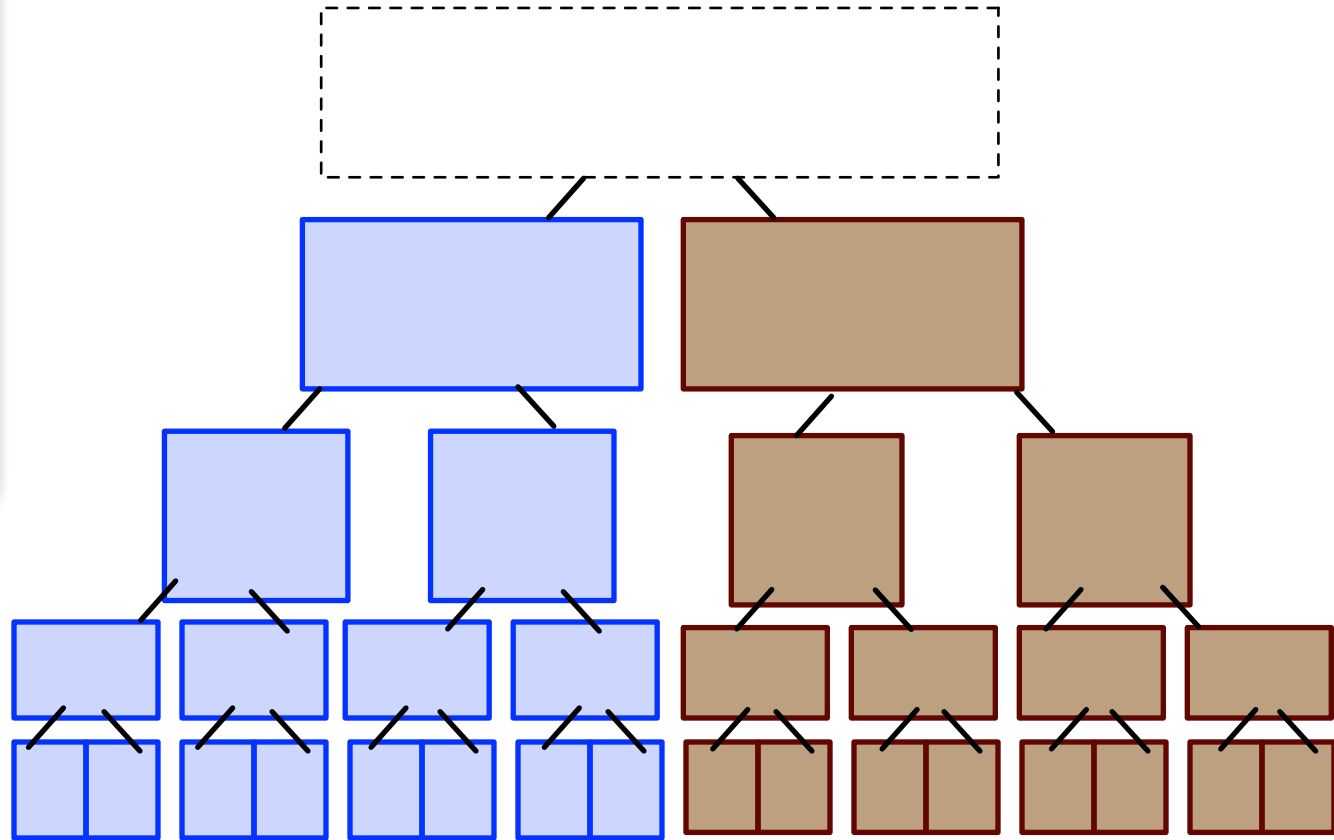
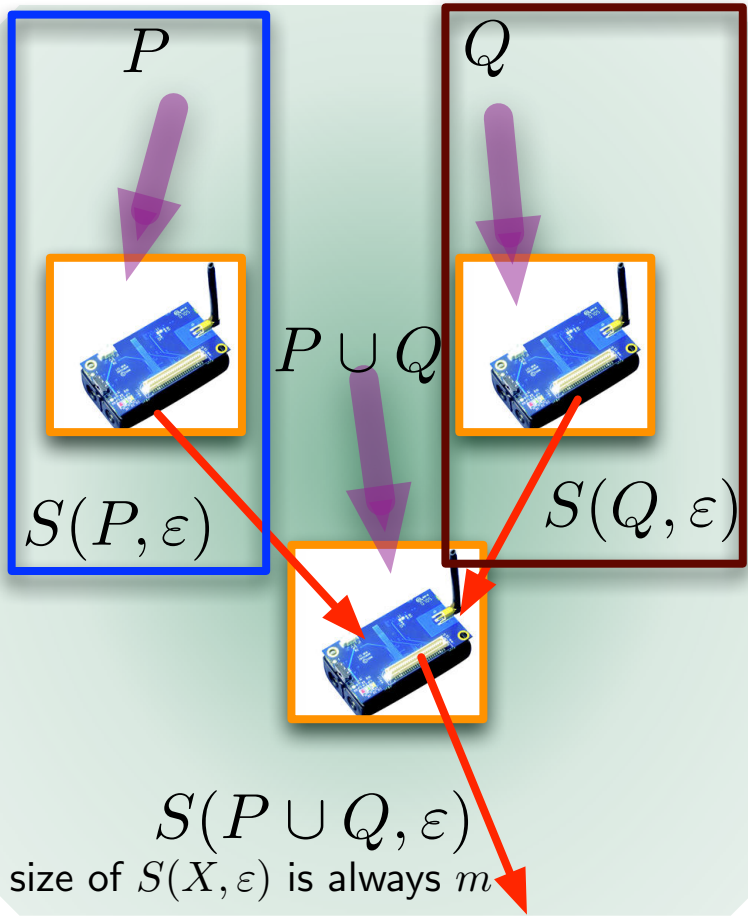
Even-Weight Merge: $(1/\epsilon) \sqrt{\log(1/\epsilon\delta)}$.



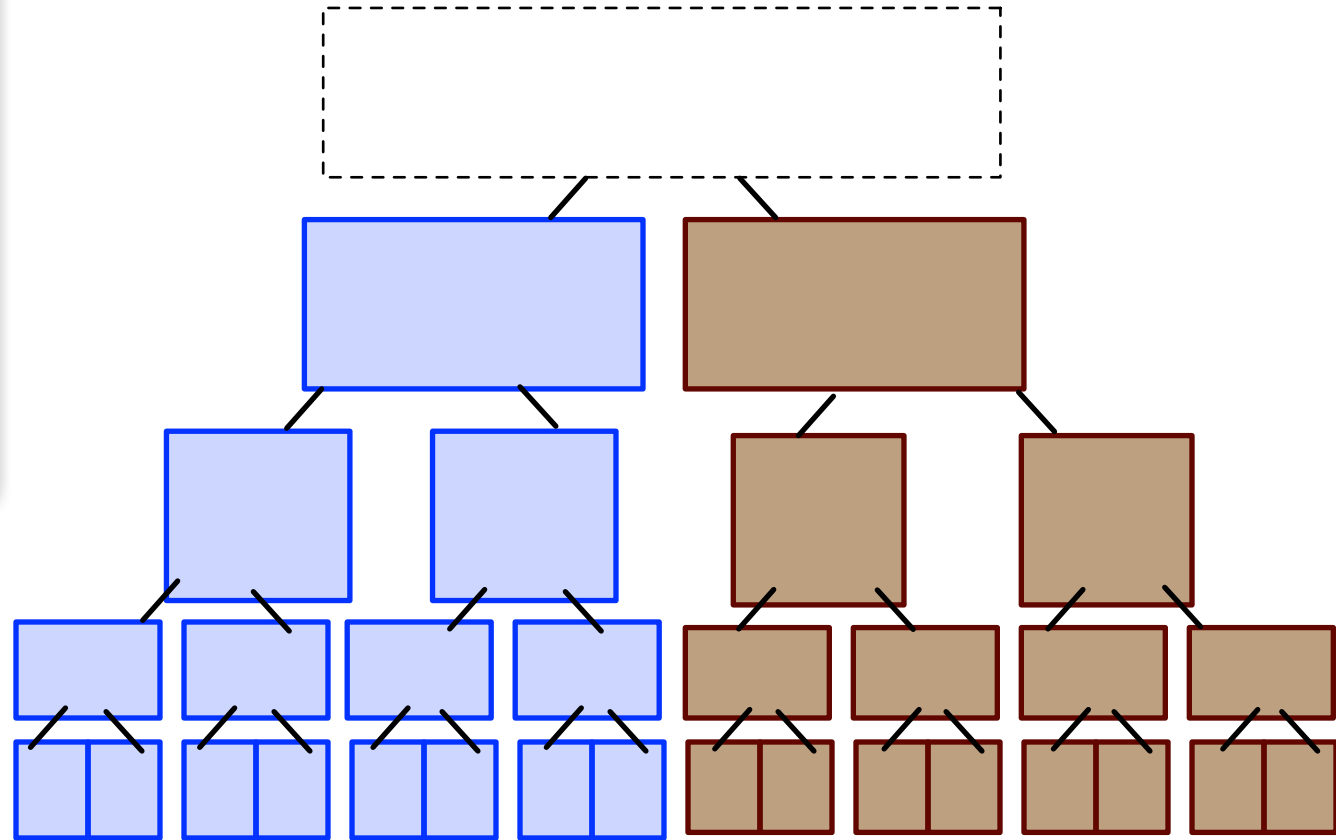
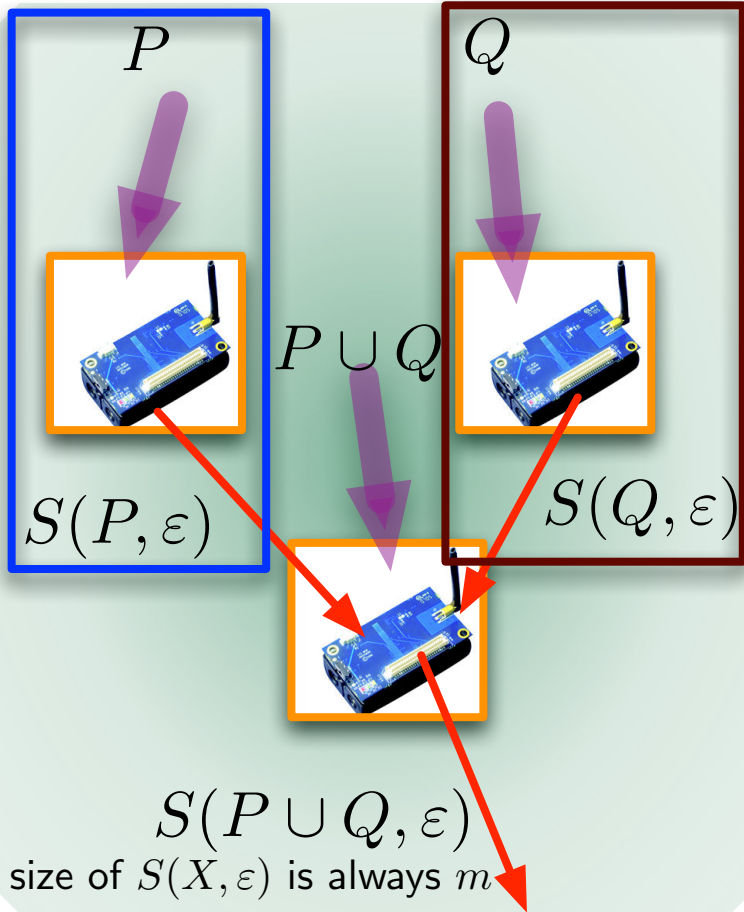
ϵ -Samples (Intervals)



ϵ -Samples (Intervals)



ε -Samples (Intervals)

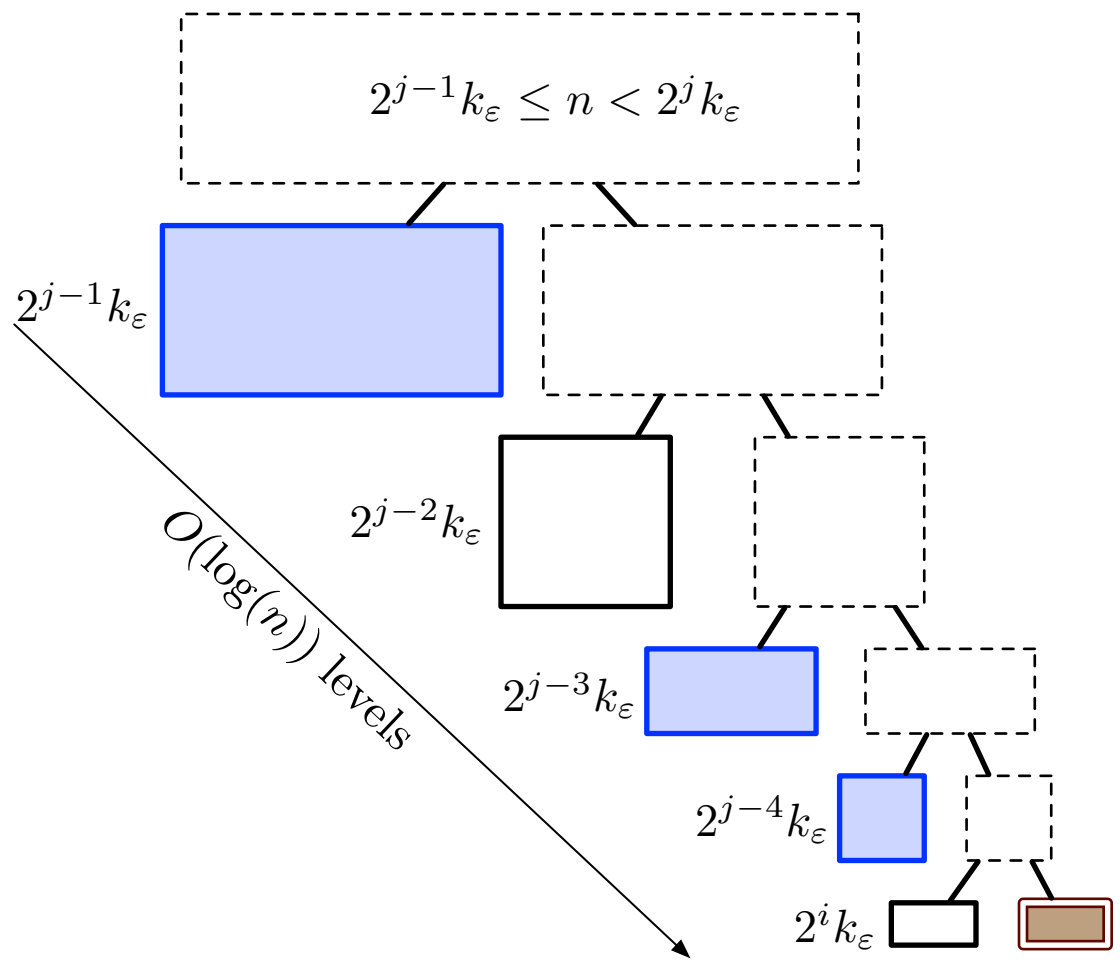
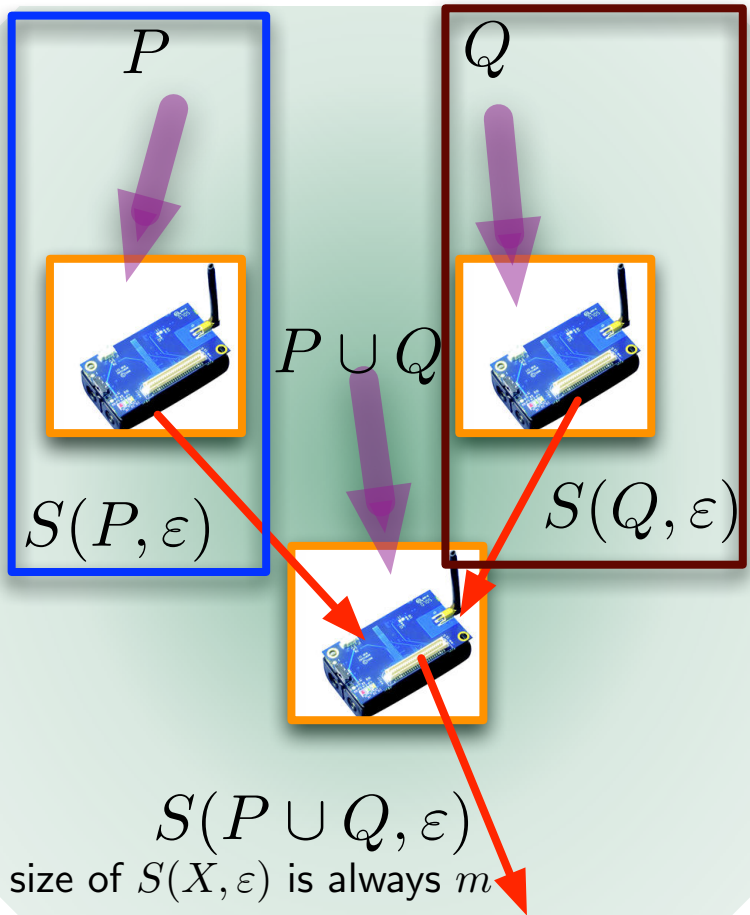


Let $E_{i,j}$ is j th merge error at level i .
 $\mathbf{E}[E_{i,j}] = 0$ and $|E_{i,j}| \leq 2^i = \Delta_i$

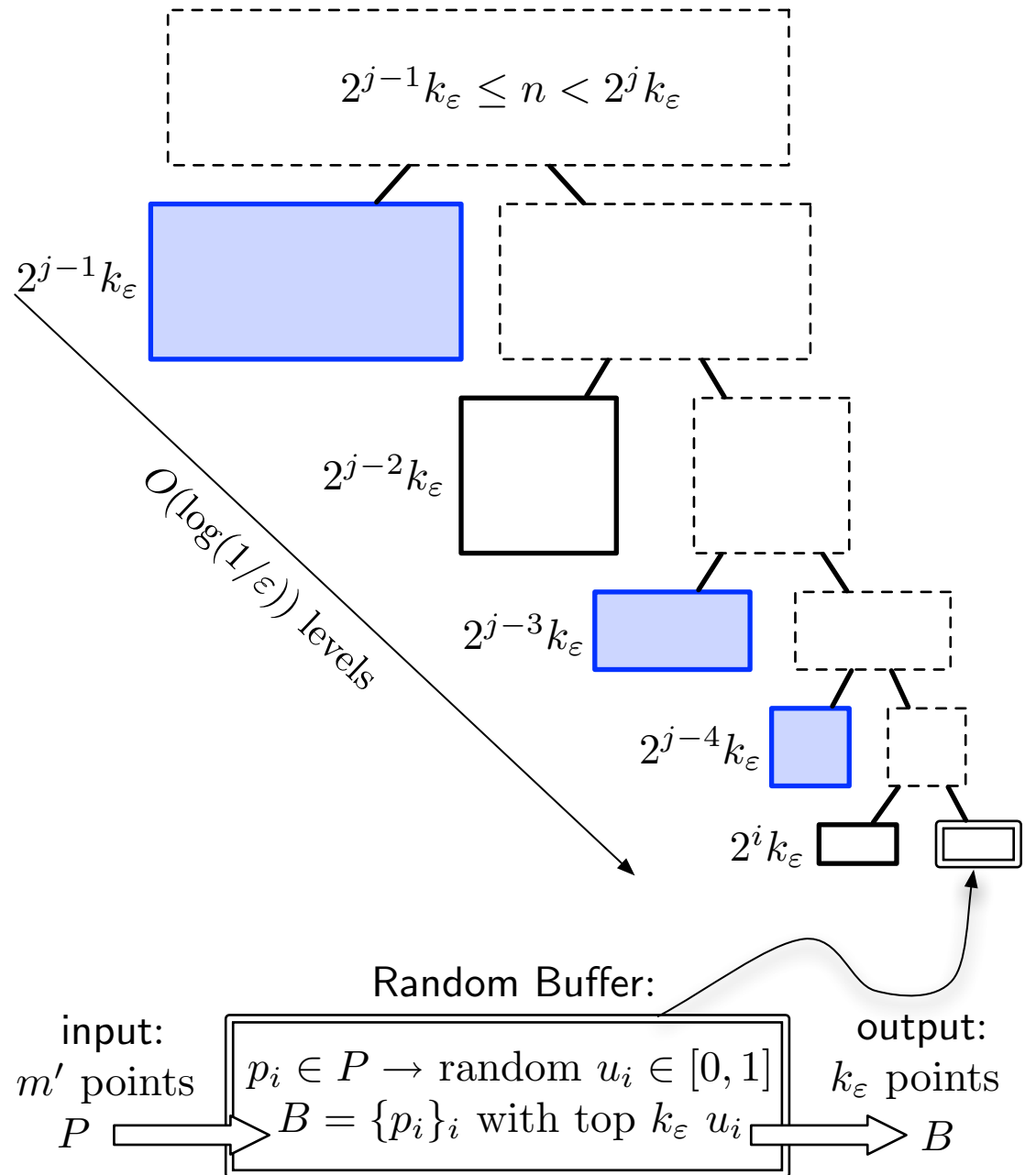
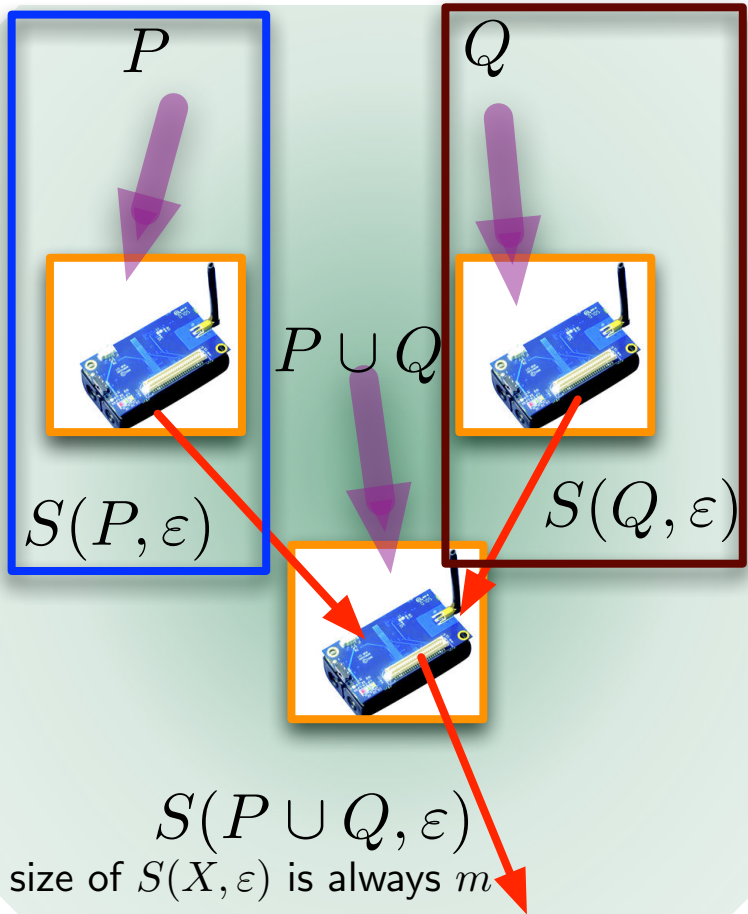
Chernoff-Hoeffding Bound:

$$\Pr[\text{ERR} > \varepsilon] \leq 2 \exp\left(\frac{-2\varepsilon^2}{\sum_i \sum_j \Delta_j^2}\right) \leq \delta$$

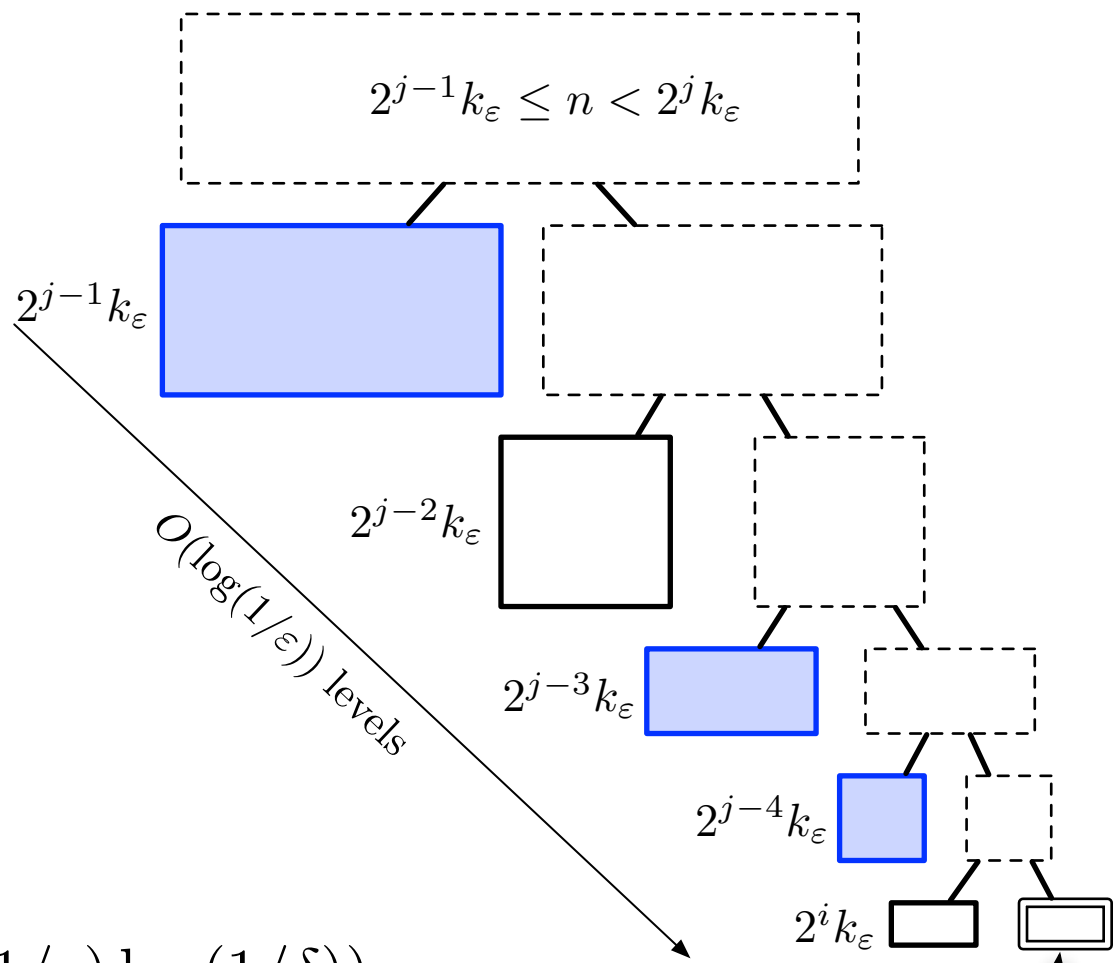
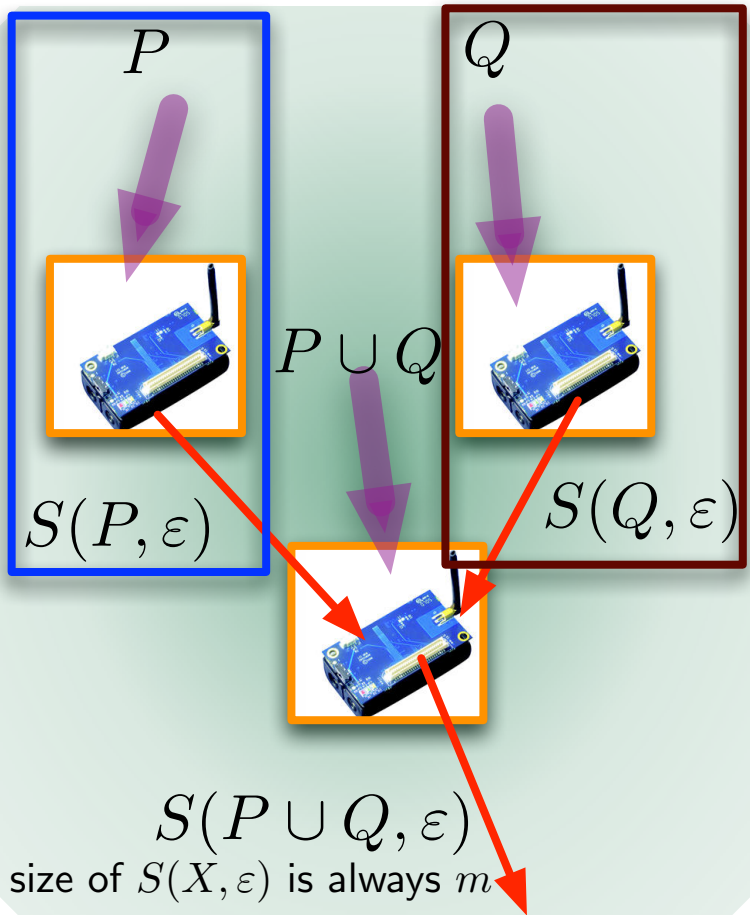
ϵ -Samples (Intervals)



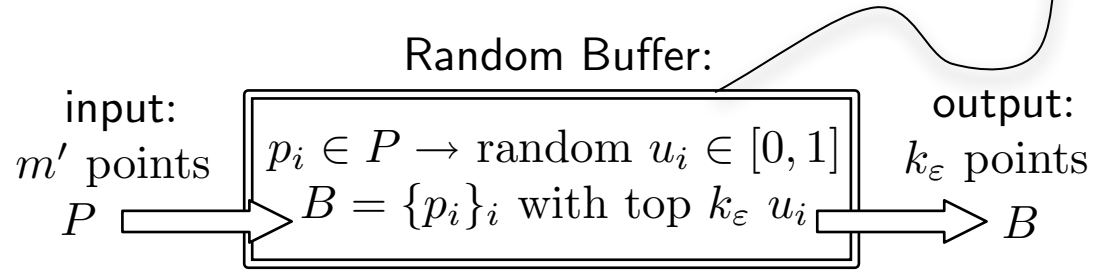
ϵ -Samples (Intervals)



ϵ -Samples (Intervals)



$$m = O\left(\frac{1}{\epsilon} \log^{1.5}\left(\frac{1}{\epsilon}\right) \log\left(\frac{1}{\delta}\right)\right)$$



Mergeable Summaries for MASSIVE Data

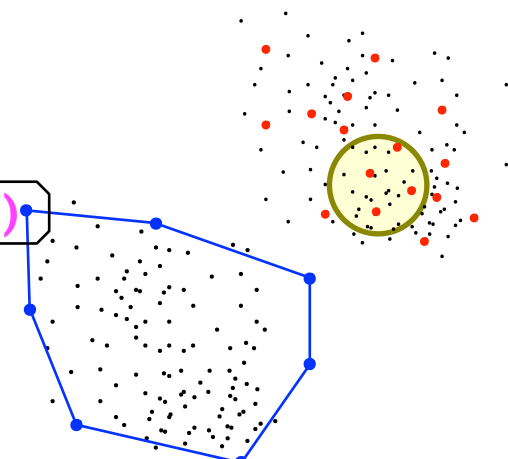
Allows approximate computation with guarantees and small space

coreset: small summary, proxy for full data set

with approx guarantees:

- ε -samples of (P, \mathcal{R}) : approx density **Mergeable**

- ε -kernel: approx convex shape **Mergeable (restricted)**



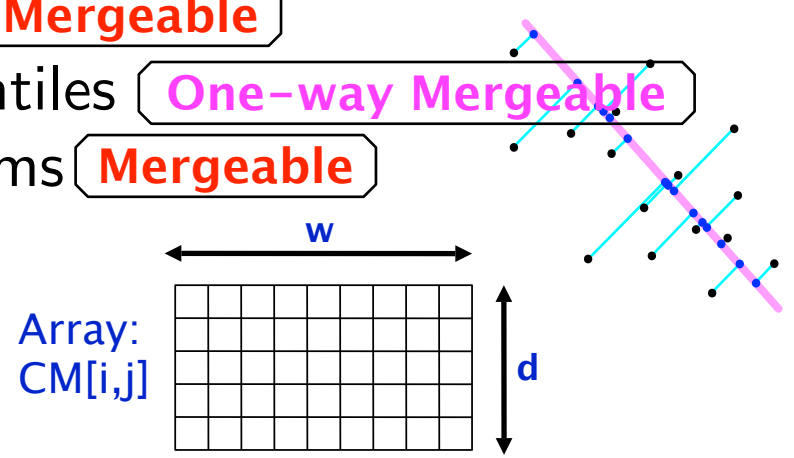
sketch: (random) (linear) combination of full data, recover functions with approx guarantees:

- Euclidean distance: Johnson-Lindenstrauss random projection **Mergeable**

- min-count sketch: approx item counts **Mergeable**

- Greenwald-Khanna sketch: approx quantiles **One-way Mergeable**

- Misra-Gries sketch: approx frequent items **Mergeable**



Open Questions

- Mergeable ε -kernels without restrictions
- Mergeable summaries for clustering
- Mergeable summaries for PCA
- Mergeable summaries for graphs [next talk]
- Lower bounds for mergeable summaries (deterministic)
- Implementation Studies

