

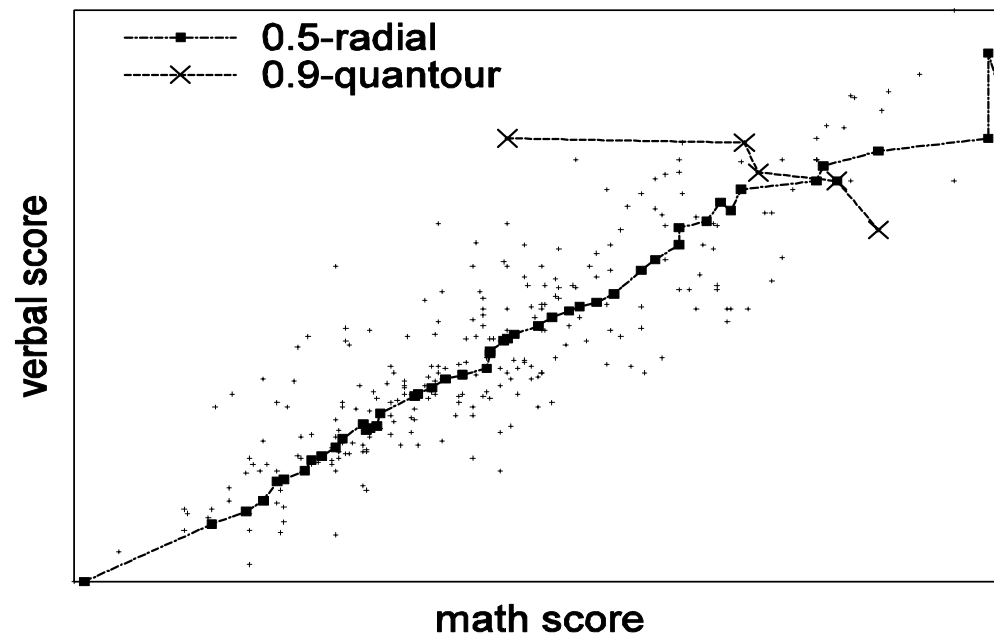
Summarizing 2D Data with Skyline-based Statistical Descriptors

Flip Korn, AT&T Labs

joint work with
Graham Cormode, Divesh Srivastava
AT&T Labs

Example: SAT Scores

- Percentiles on each variable not sufficient
- How to compare values across dimensions?
- Dominant students, balanced students



Example: Flow Size Distribution

- IP flow summarizes #pkts, #bytes
- Joint distribution indicates behavior
- Track changes in bandwidth, application, etc.
- Detect anomalies (DoS attacks, BGP flaps)

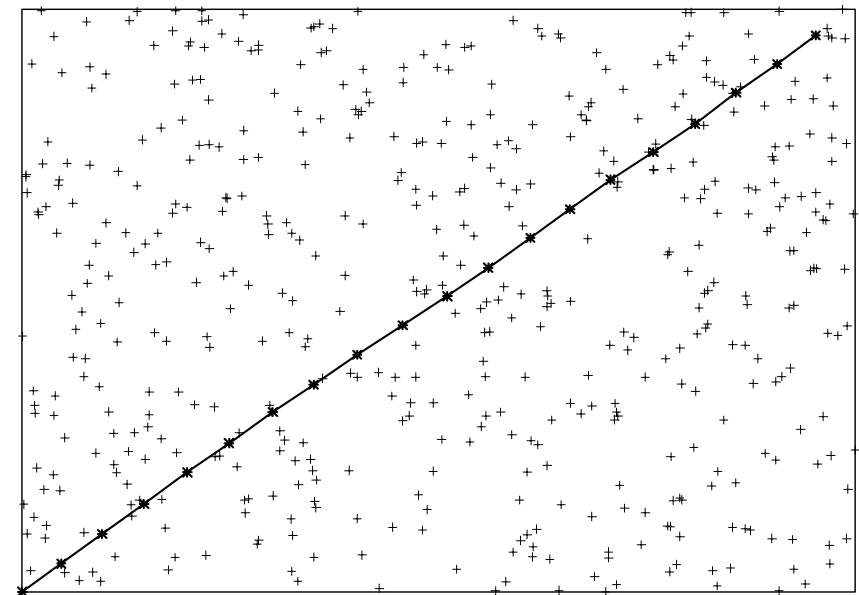
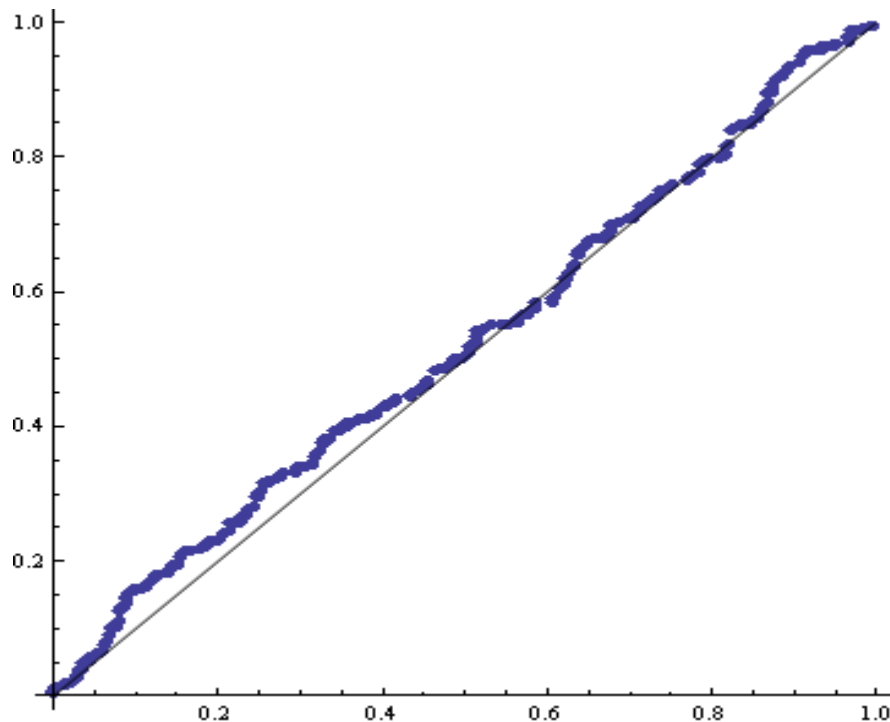
	Many pkts	Few pkts
Many bytes	ftp	http
Few bytes	ping	dhcp

Desiderata

- Quantile: item with rank $\Gamma \notin N$ (eg, median)
- But no total (rank) ordering in 2D
- Goal: capture joint distribution
 - Dominance: Pareto optimality
 - Skew: compare trade-offs betw variables
- Robustness is crucial
 - based on *rank*s, not values

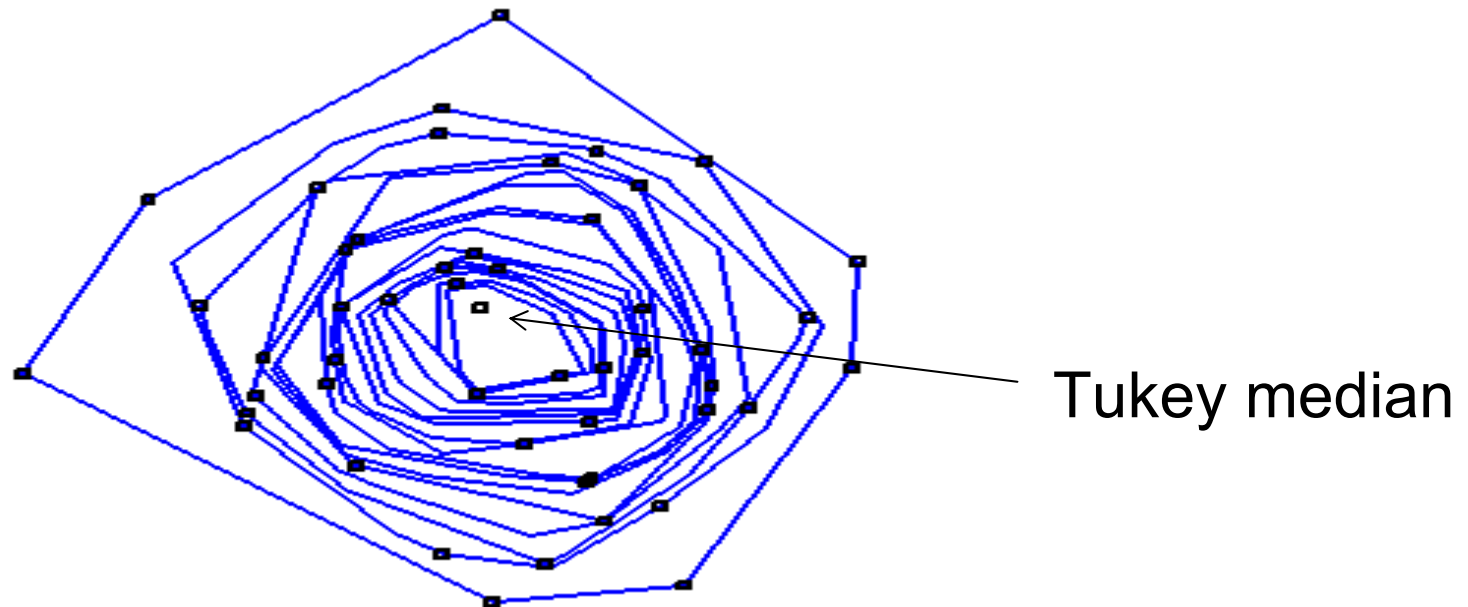
QQ-Plots

- Quantiles on each dimension
- Doesn't consider joint distribution
- May return points not in data set



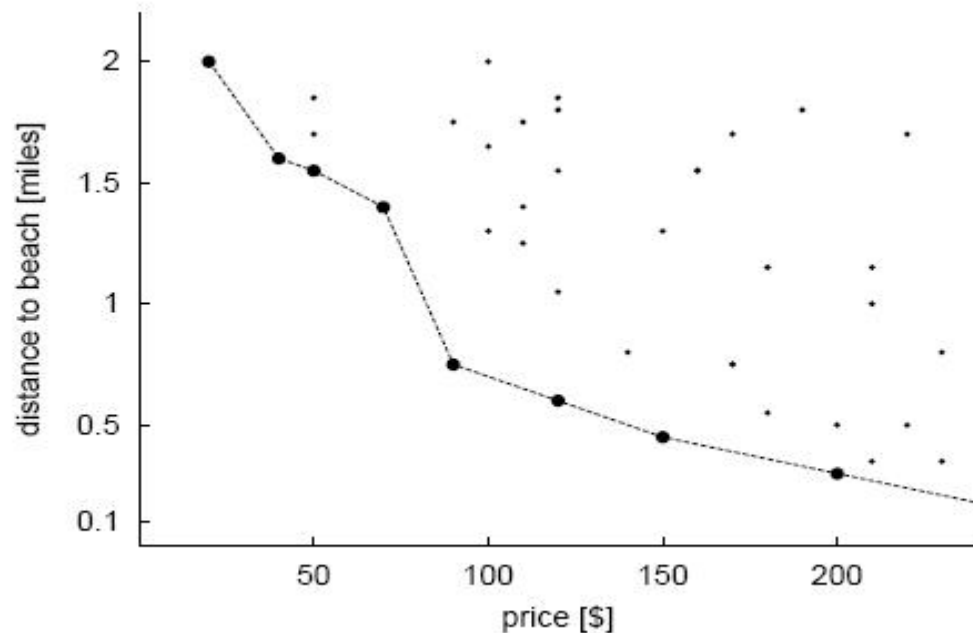
Depth Contours

- Recursively compute convex hull
- Dot-product dominance
- Arbitrary #layers, #points per layer



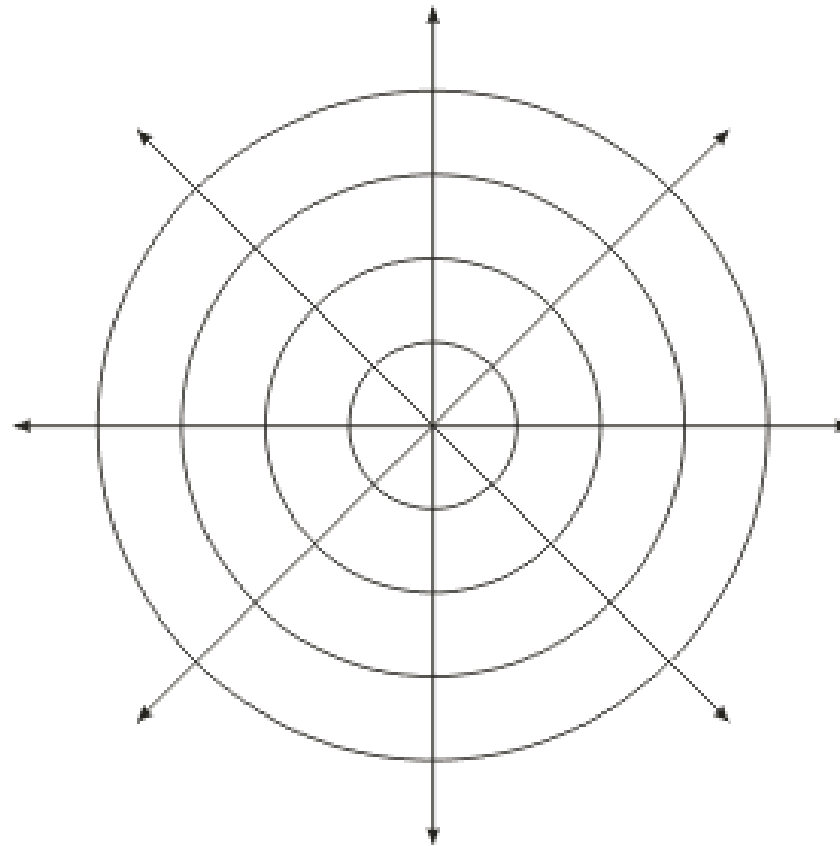
Skylines

- Points not dominated by a point (k points)
- Can be of arbitrary size
- Additional criteria, but not based on distribution



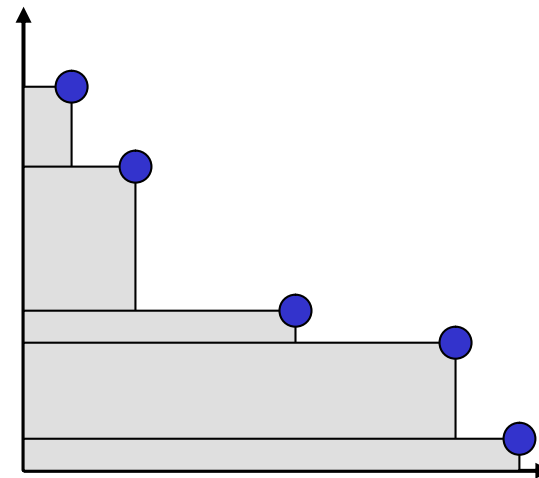
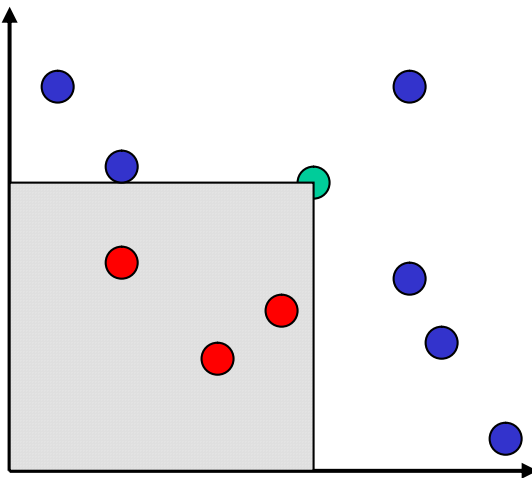
Our idea

- Simple and natural: inspired by polar coords!



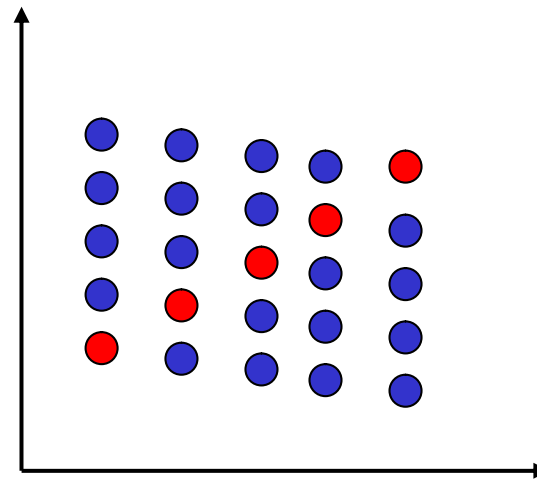
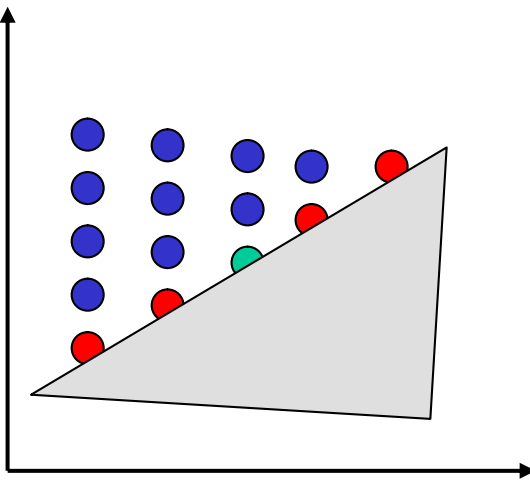
ϕ -Quantour Definition

- p **dominates** q iff $(p_x \geq q_x)$ and $(p_y \geq q_y)$
- $P \phi =$ points dominating $\leq \Gamma \phi N$ points
- ϕ -quantour: skyline of $P \phi$



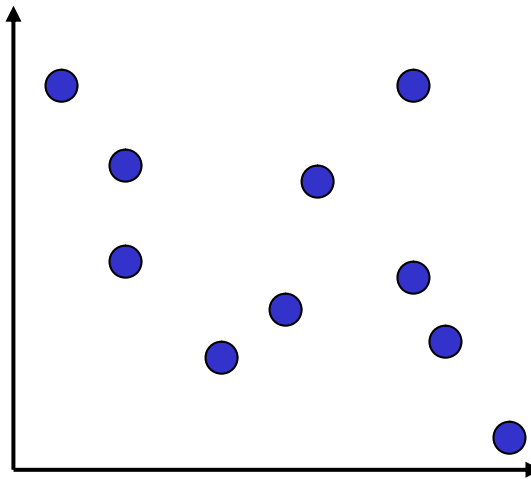
α -Radial Definition

- **skew** = $\text{ranky}(p) / [\text{xrank}(p) + \text{yrank}(p)]$
- P_α = points with $\text{skew} \leq \Gamma_\alpha N$ points
- α -radial: skyline of P_α



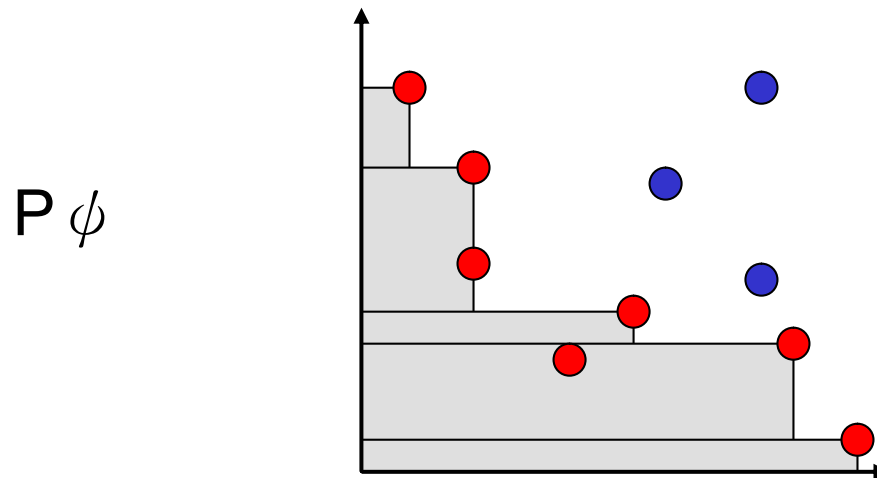
(α, ϕ) -Quantile Definition

- ϕ -skyline of α -skyline of $P \cap P_\alpha$



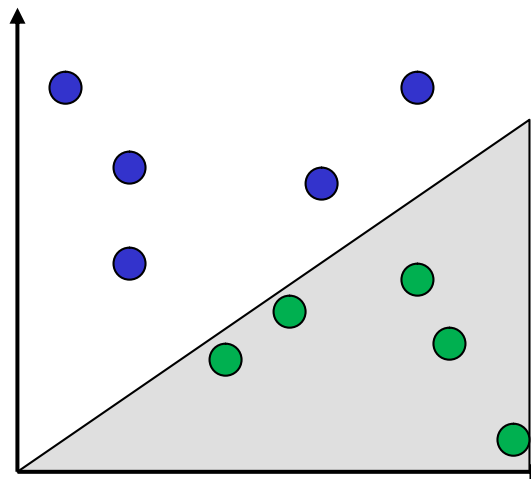
(α, ϕ) -Quantile Definition

- ϕ -skyline of α -skyline of $P_\phi \cap P_\alpha$



(α, ϕ) -Quantile Definition

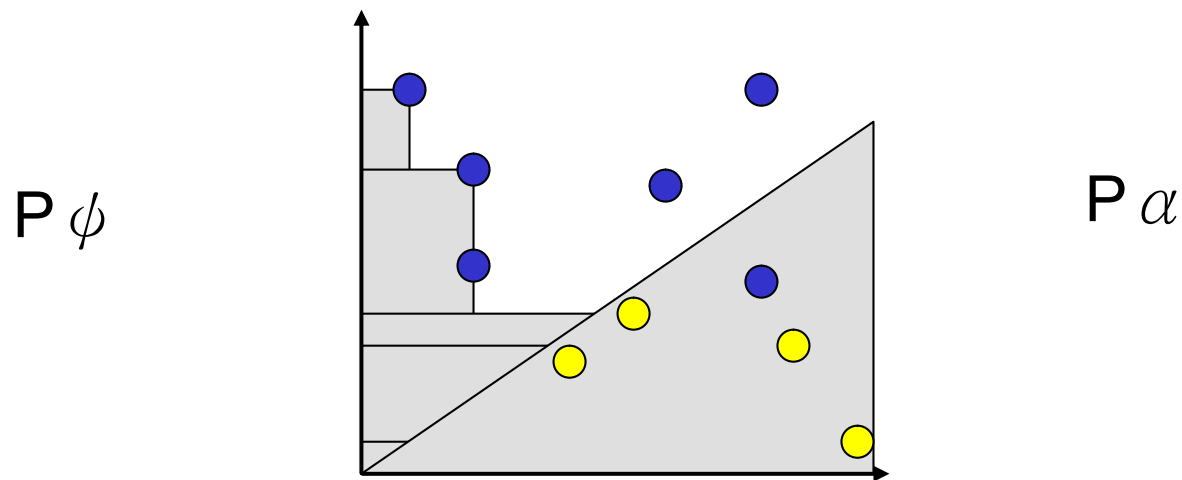
- ϕ -skyline of α -skyline of $P_\phi \cap P_\alpha$



P_α

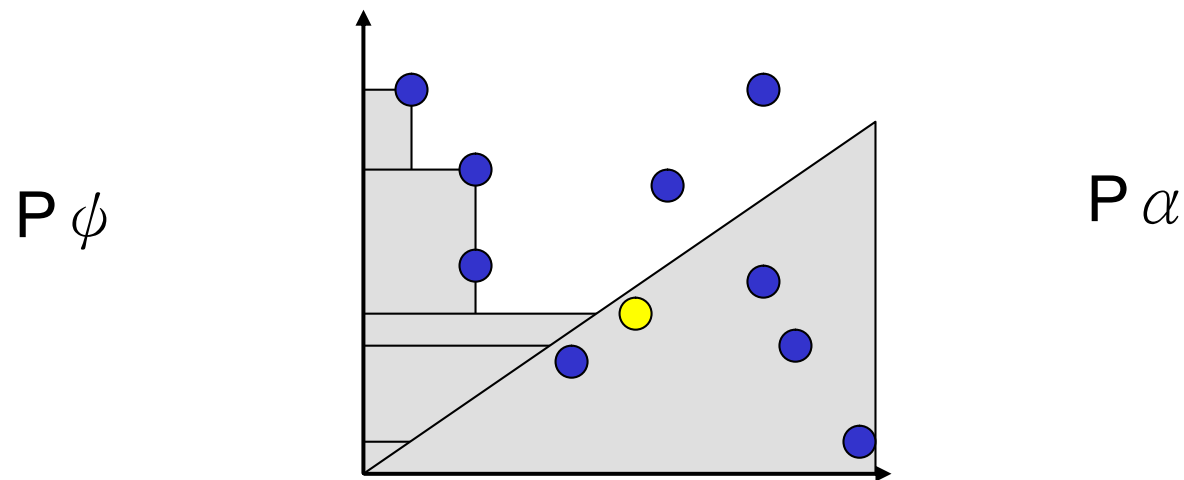
(α, ϕ) -Quantile Definition

- ϕ -skyline of α -skyline of $P_\phi \cap P_\alpha$



(α, ϕ) -Quantile Definition

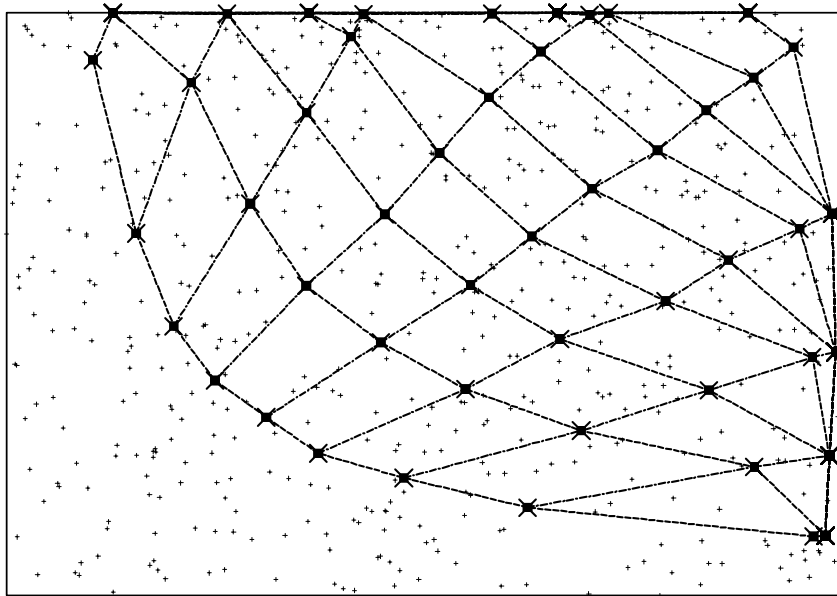
- ϕ -skyline of α -skyline of $P_\phi \cap P_\alpha$



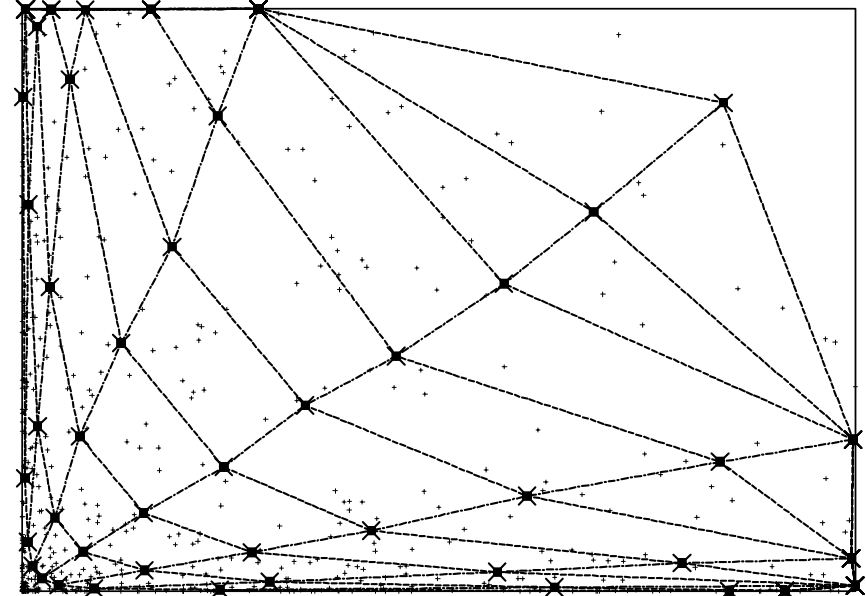
Properties of (α, ϕ) -Quantiles

- Order of skylines on $P_\phi \cap P_\alpha$ doesn't matter
- p lies on α -radial or ϕ -quantour or both
- Uniqueness
 - Unique point p always found
 - $p = q$ iff $\phi(p) = \phi(q)$ and $\alpha(p) = \alpha(q)$
- Collapses to 1D quantiles for points (x,x)

Illustration

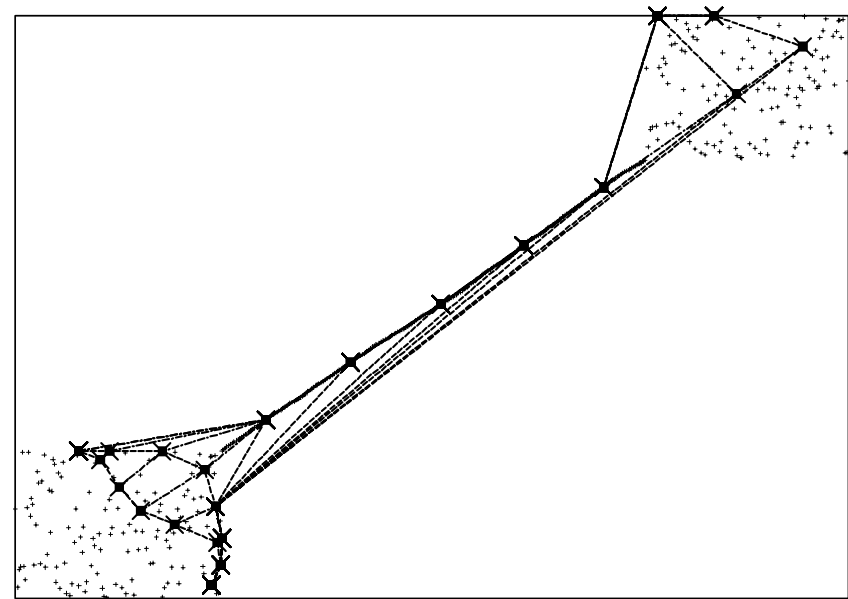
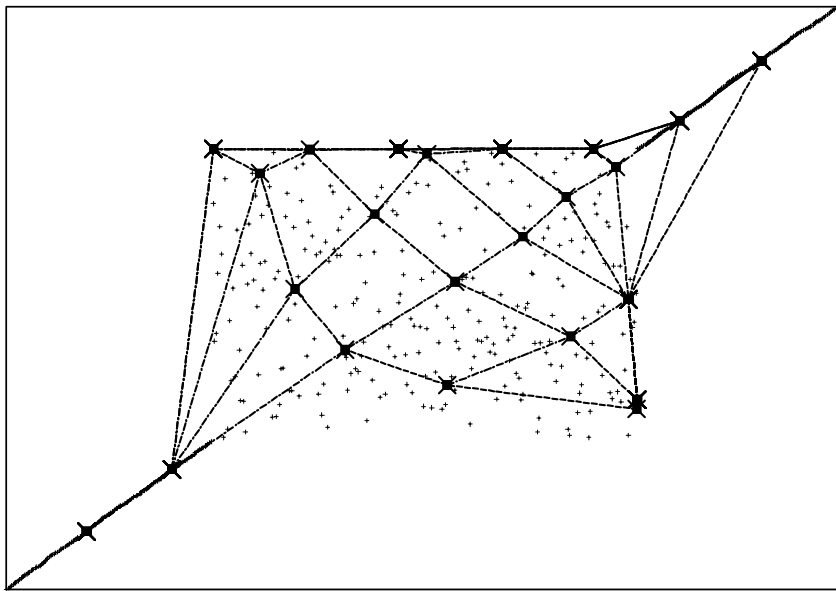


Uniform x Uniform

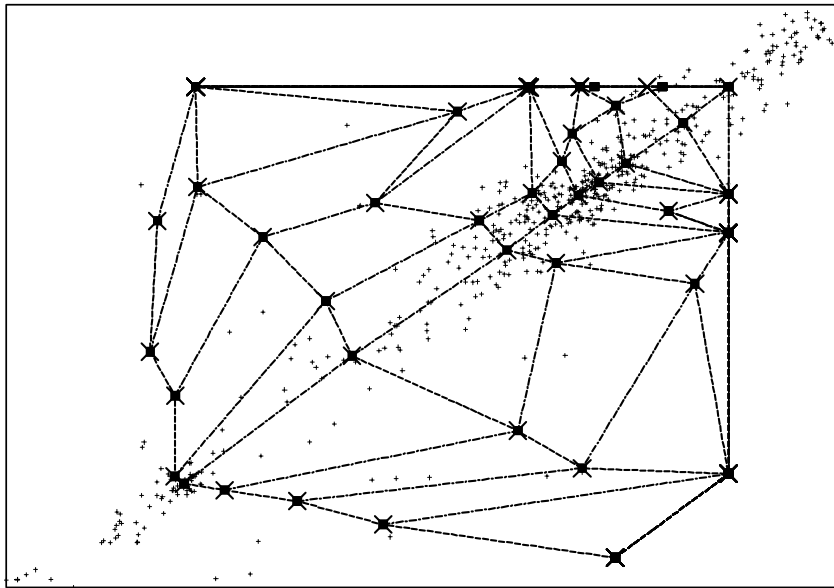


Zipf x Zipf

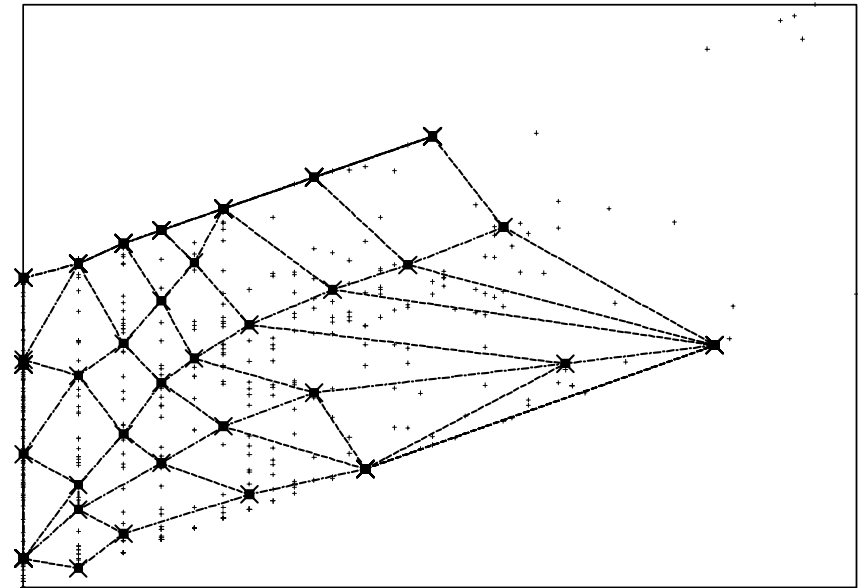
Illustration



Illustration



SNMP data



Flow data (log-log)

Exact Algorithm

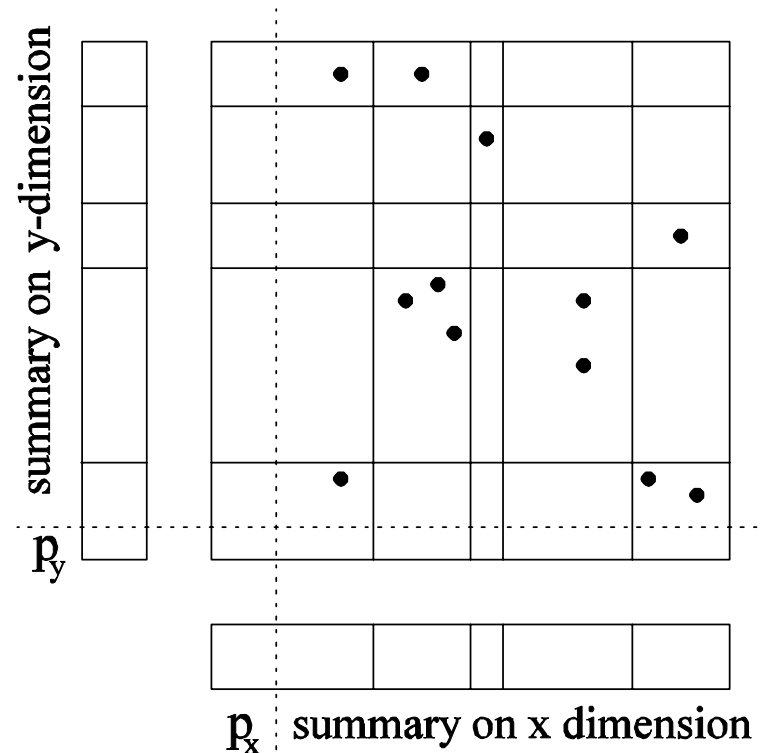
- Given (α, ϕ) find (i_α, j_ϕ) -quantiles
- Pre-compute $xrank(p)$, $yrank(p)$, $rank(p)$
- Scan each point p and compute
 - $i = \Gamma_\alpha(p) / \alpha$, $j = \Gamma_\phi(p) / \phi$
 - Maintain point with max- y (max- x if tie) at (i,j)
 - Compute prefix max along rows/cols of (i,j)
- Complexity is $O(N \log N)$

Streaming Algorithms

- Rank computation: given p , find $(\alpha(p), \phi(p))$
- Primitives: GK-digest, QD-digest, BQ-digest
- Three approaches
 - Cross-product
 - Deferred-merge
 - Eager-merge
- Two based on existing algorithms, one is novel
- Trade-offs (space, time, skew, etc)

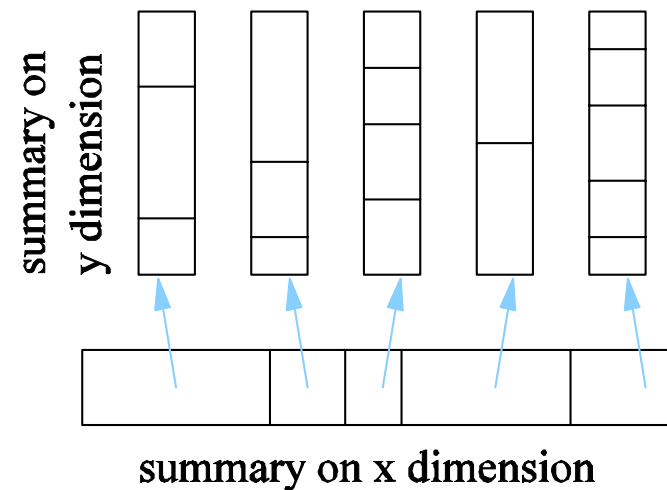
Cross-Product Approach

- 1D digest on each axis
- Cross-product grid
- Insert row + col
- Merge whole rows/cols
- GK x GK
 - $O(1/\epsilon^2 \log^2(\epsilon N))$ space
 - $O(1/\epsilon \log \epsilon N)$ update



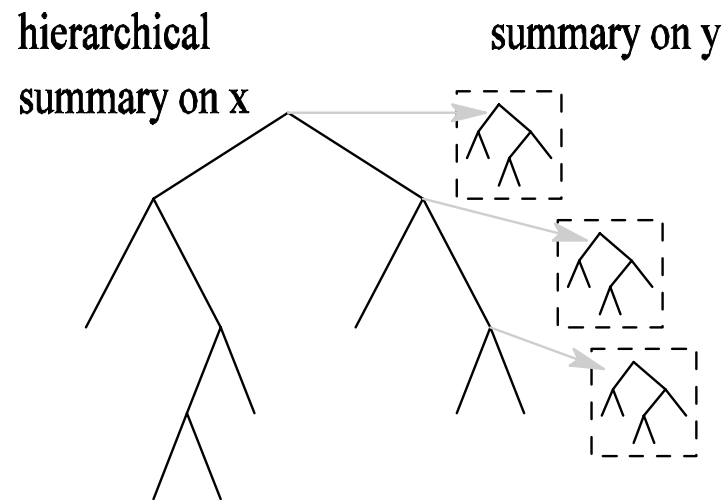
Deferred-Merge Approach

- Cascaded 1D digests
- QK x QD
 - $O(1/\epsilon^2 \log(\epsilon N) \log U)$ sp
 - $O(\log(1/\epsilon \log U))$ update



Eager-Merge Approach

- Cascaded 1D Digests
- Multiple insertions
- No merging
- QD x QD
 - $O(1/\epsilon \log^3 U)$ space
 - $O(\log U \log \log U)$ update

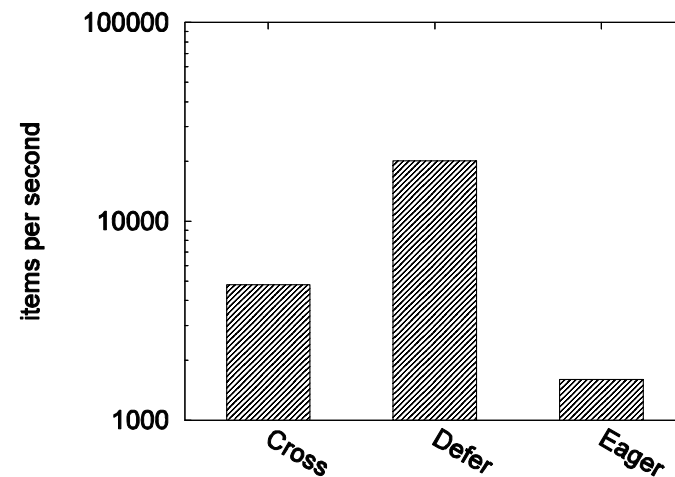
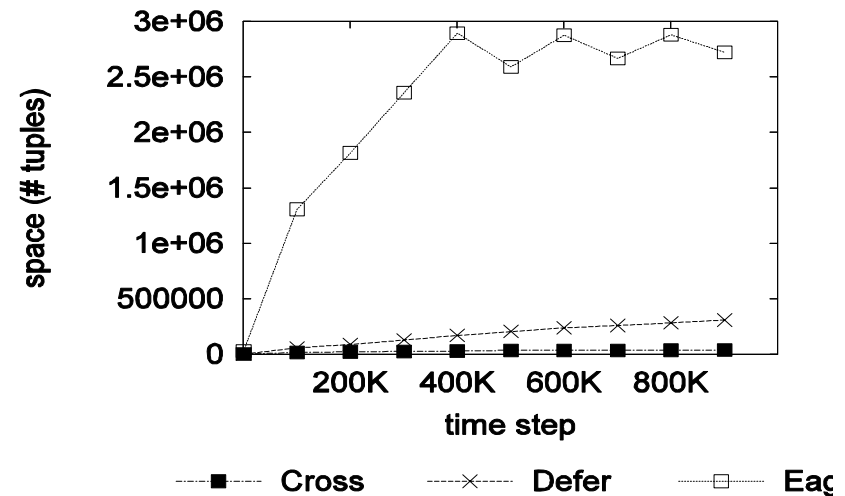


Summary of Results

	Space	Update
Cross-Product (GK x GK)	$1/\varepsilon^2 \log^2(\varepsilon N)$	$1/\varepsilon \log(\varepsilon N)$
Deferred-Merge (GK x QD)	$1/\varepsilon^2 \log(\varepsilon N) \log U$	$\log(1/\varepsilon \log U)$
Eager-Merge (QD x QD)	$1/\varepsilon \log^3 U$	$\log U \log \log U$

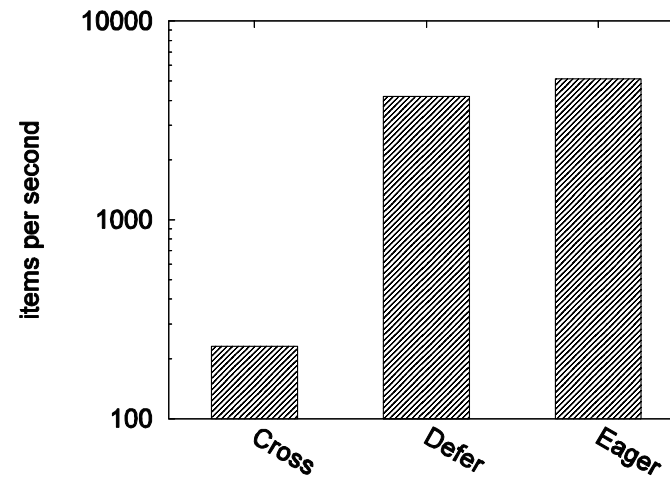
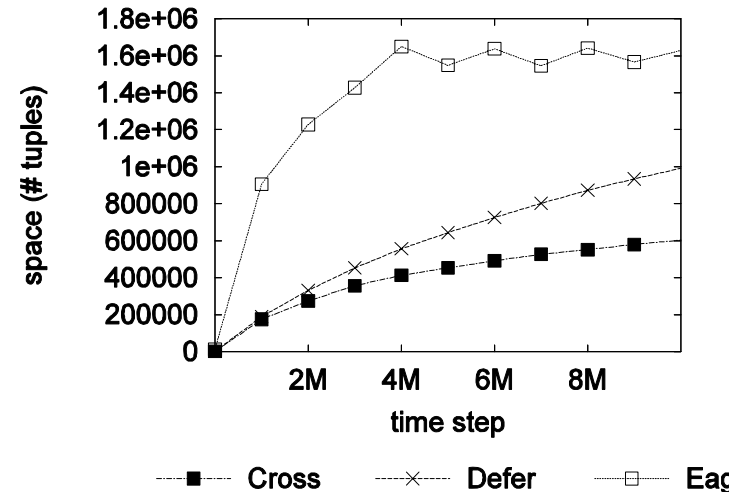
Experiments

- Flow data
- $\varepsilon = 0.01$
- Space usage
 - Eager worst
- Throughput
 - Eager worst, defer best



Experiments (cont'd)

- Flow data
- $\varepsilon = 0.001$
- Space usage
 - Eager worst
- Throughput
 - Cross-product worst



Conclusions

- Proposed (α, ϕ) -quantiles
 - Capture joint distribution (skew, dominance)
 - Natural definition, declarative semantics
 - Yield single points from the data
- Study of efficient streaming algorithms
 - Deferred-merge: novel and fast
- Future work
 - analysis for selection problem in streams

The End

Definitions

- p dominates q iff $(p_x > q_x)$ and $(p_y > q_y)$
- P_ϕ = points dominated by $< \phi N$ points
- ϕ -quantour: skyline of P_ϕ
- $\text{skew} = \text{rank}_y(p) / [\text{rank}_x(p) + \text{rank}_y(p)]$
- P_α = points with $\text{skew} < \alpha N$ points
- α -radial: skyline of P_α
- (α, p) -quantile: skyline of $P_p \wedge P_\alpha$