

Online Generation of Locality Sensitive Hash Signatures



BENJAMIN VAN DURME & ASHWIN LALL



human language technology
center of excellence

JOHNS HOPKINS
U N I V E R S I T Y

**Georgia
Tech**



Data Overload



- Our access to data is growing fast

Data Overload



- Our access to data is growing faster than our ability to process it

Data Overload



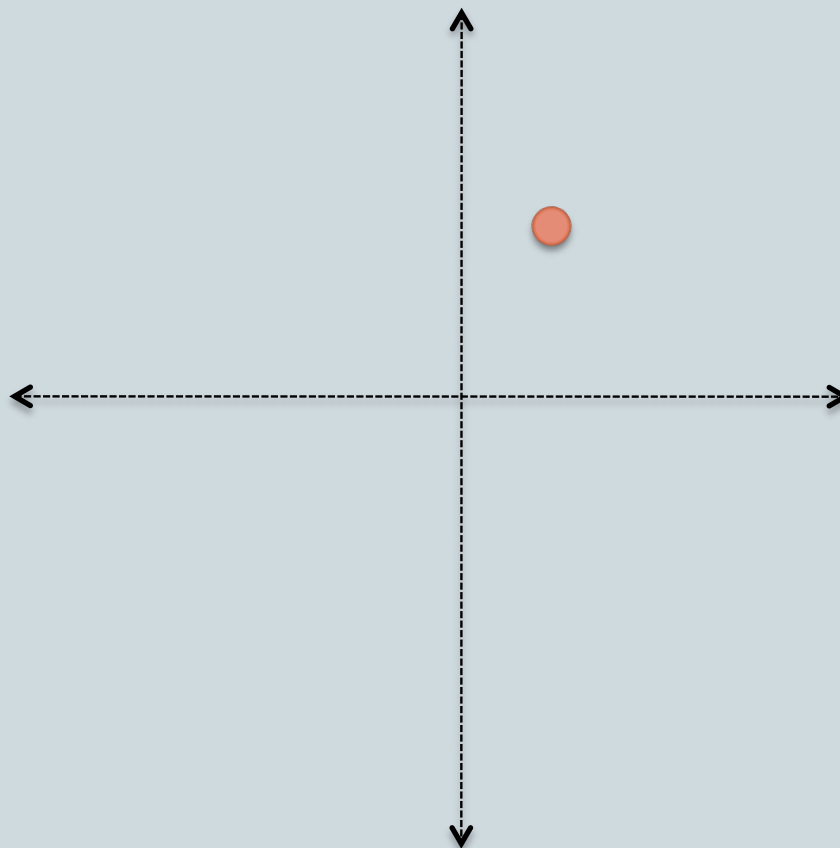
- Our access to data is growing faster than our ability to process it
- Complementary solutions:
 - Distributed environments (e.g., MapReduce)
 - Streaming / Randomized Algorithms

Locality Sensitive Hashing

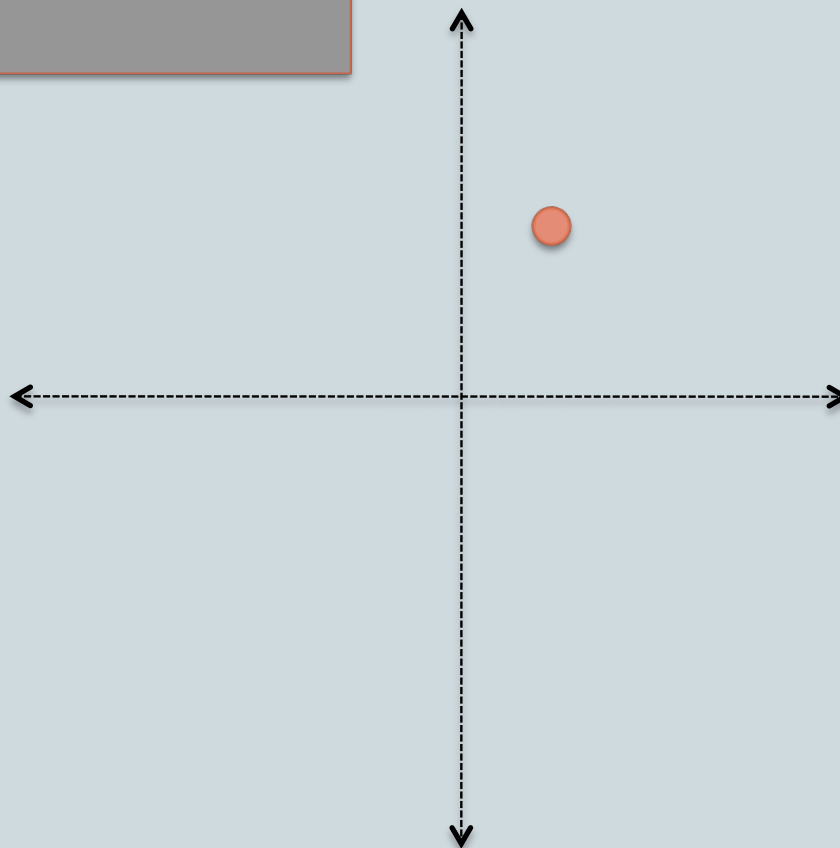
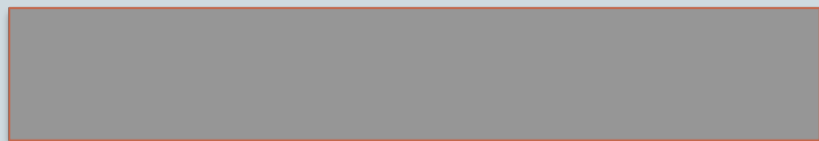


- Goal: fast comparison between points in very high dimensional space
- Indyk & Motwani ('98) => Charikar ('02)
 - Randomly project points to low dimensional *bit signatures* such that **cosine distance** is approximately preserved
- Example Applications in HLT
 - Noun clustering [Ravichandran *et al* '05]
 - Topic Detection and Tracking (TDT) [Petrovic *et al* '10]

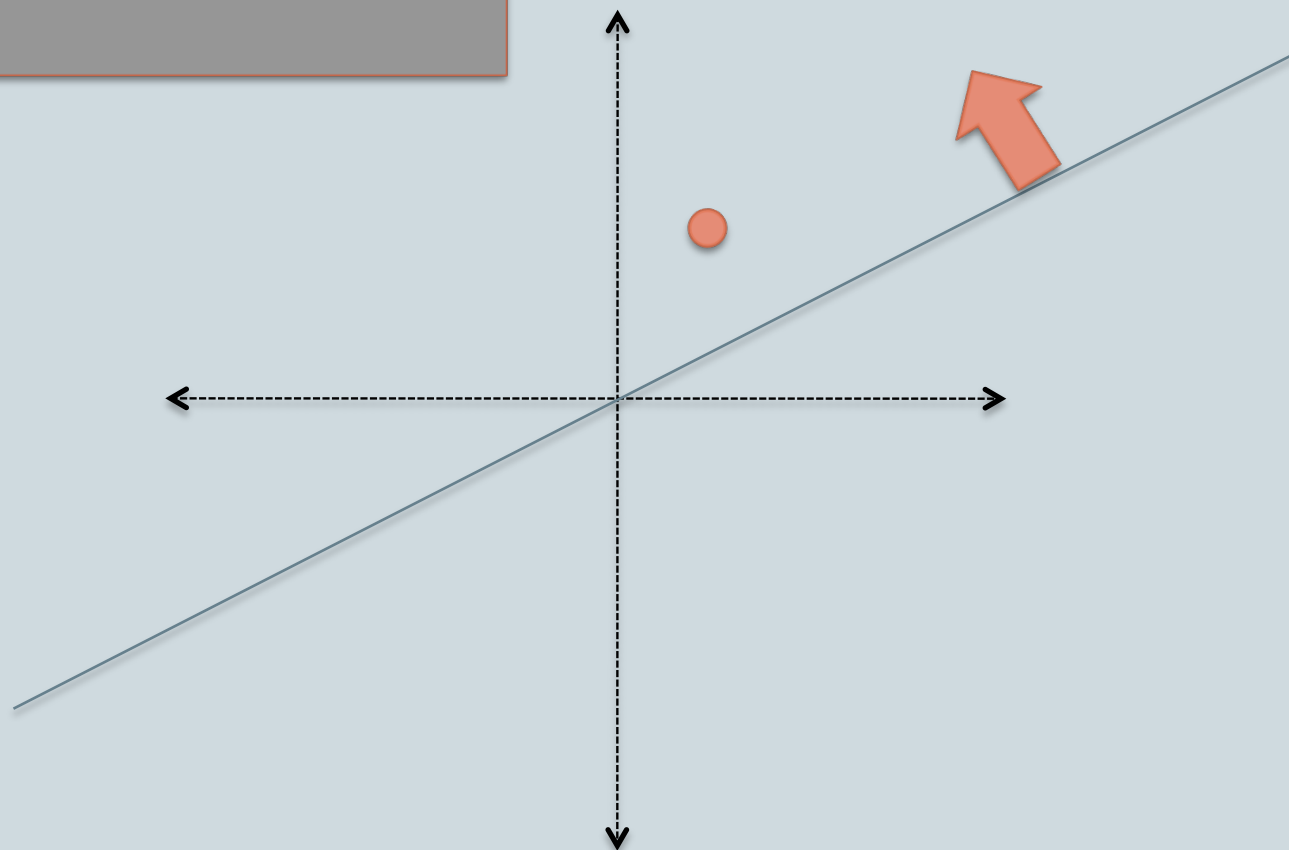
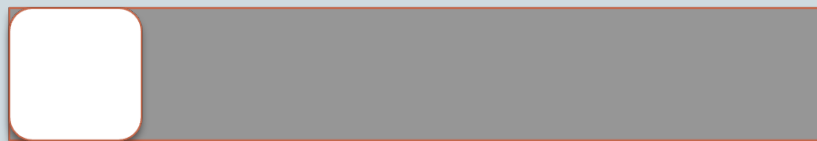
Locality Sensitive Hashing



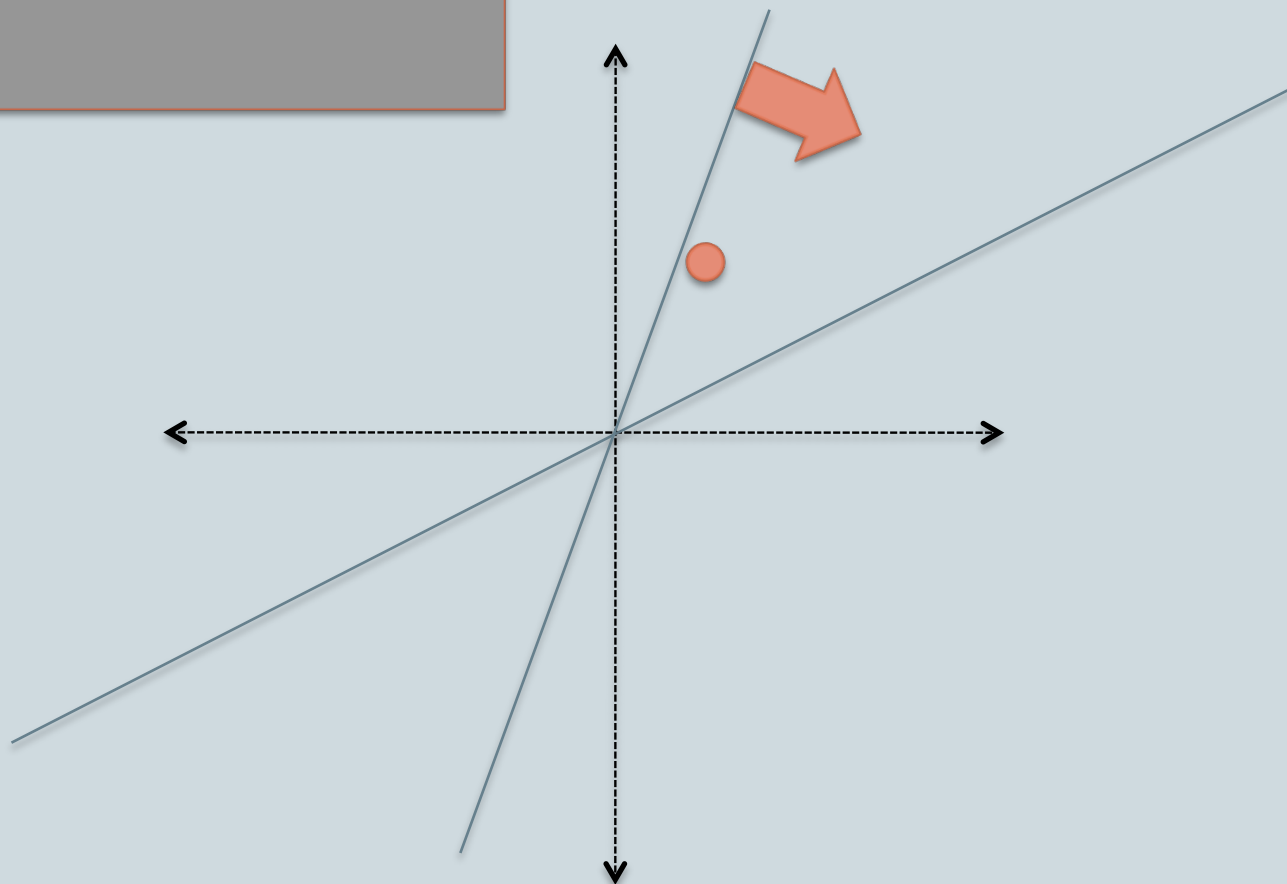
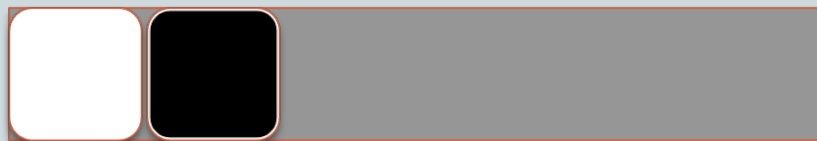
Locality Sensitive Hashing



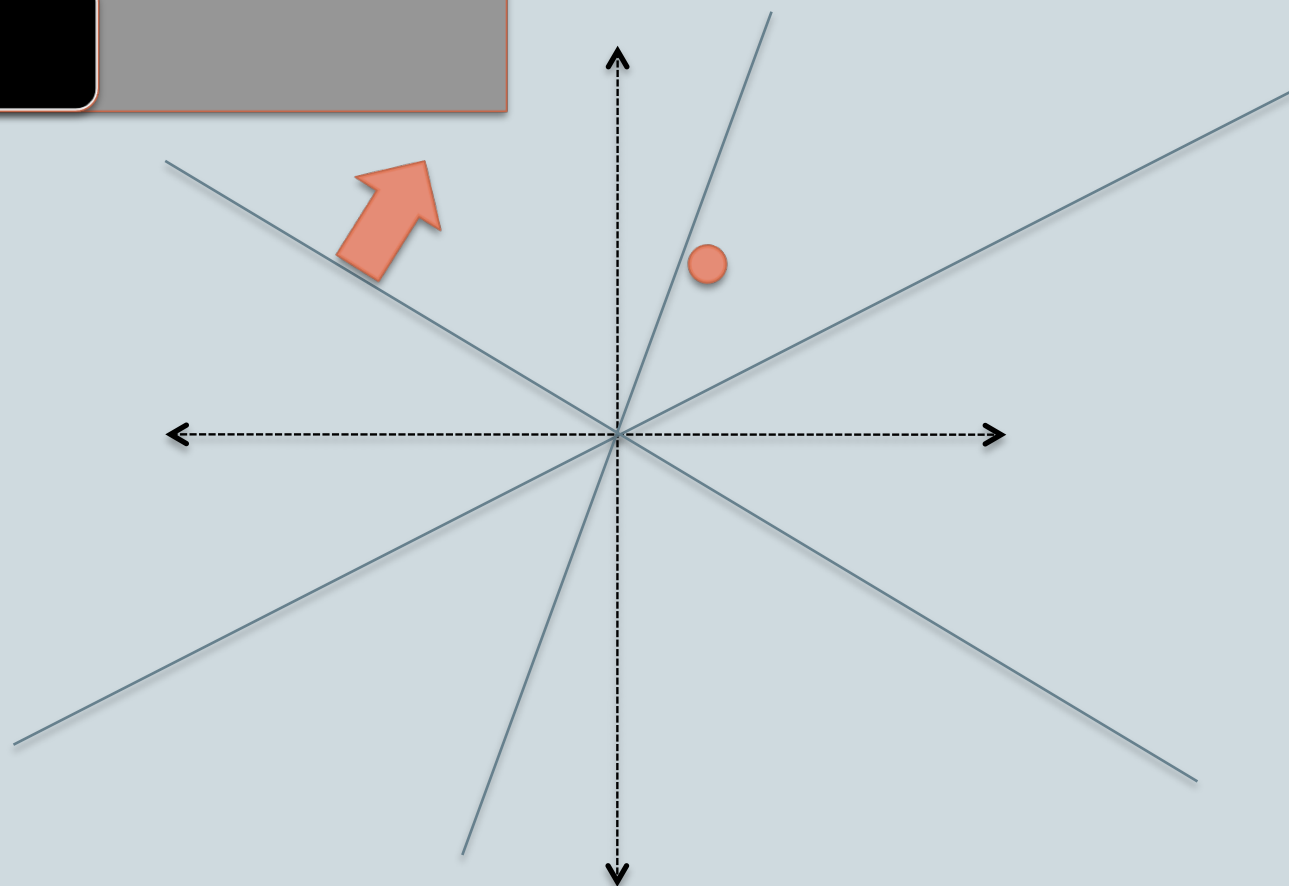
Locality Sensitive Hashing



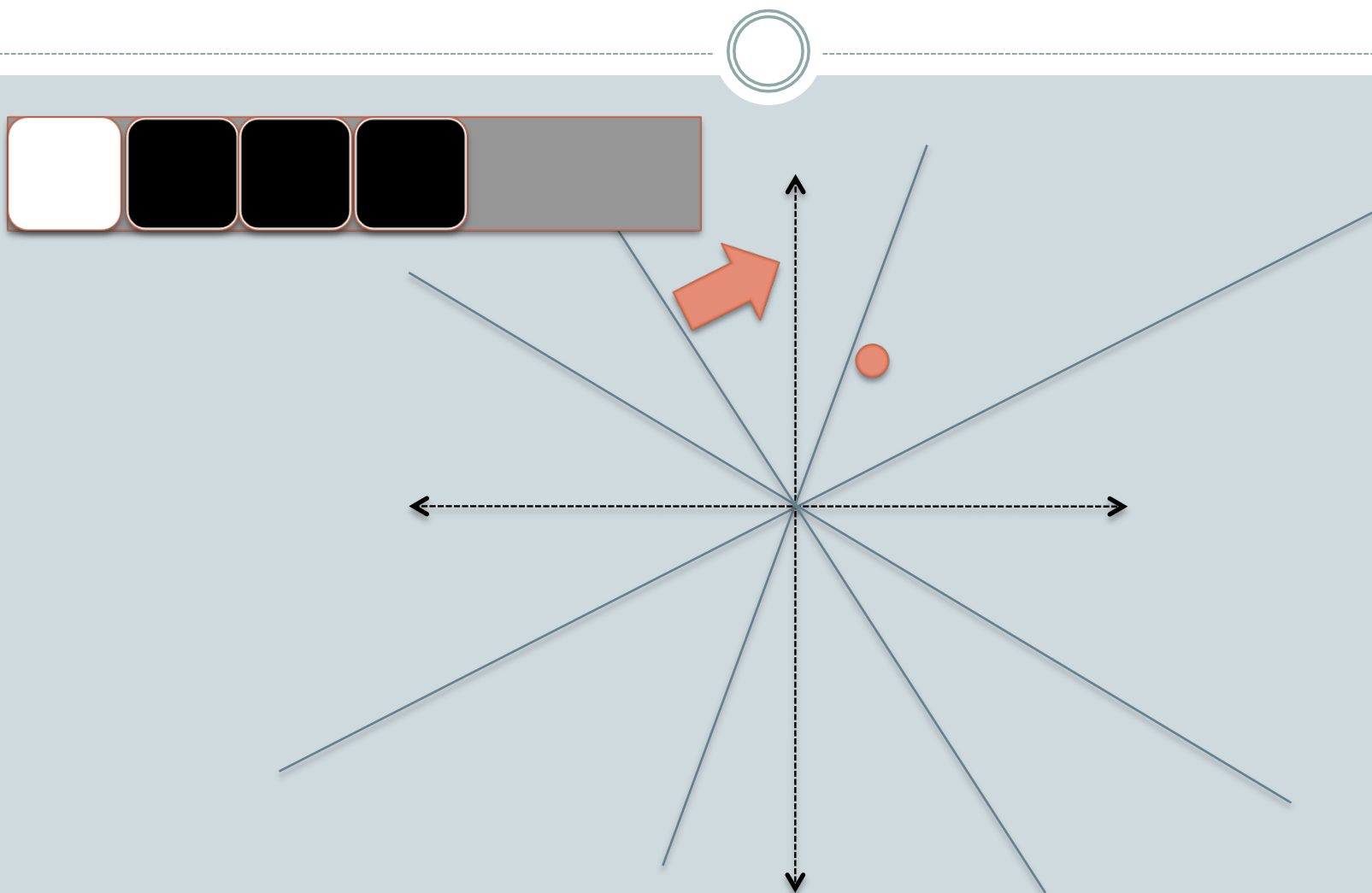
Locality Sensitive Hashing



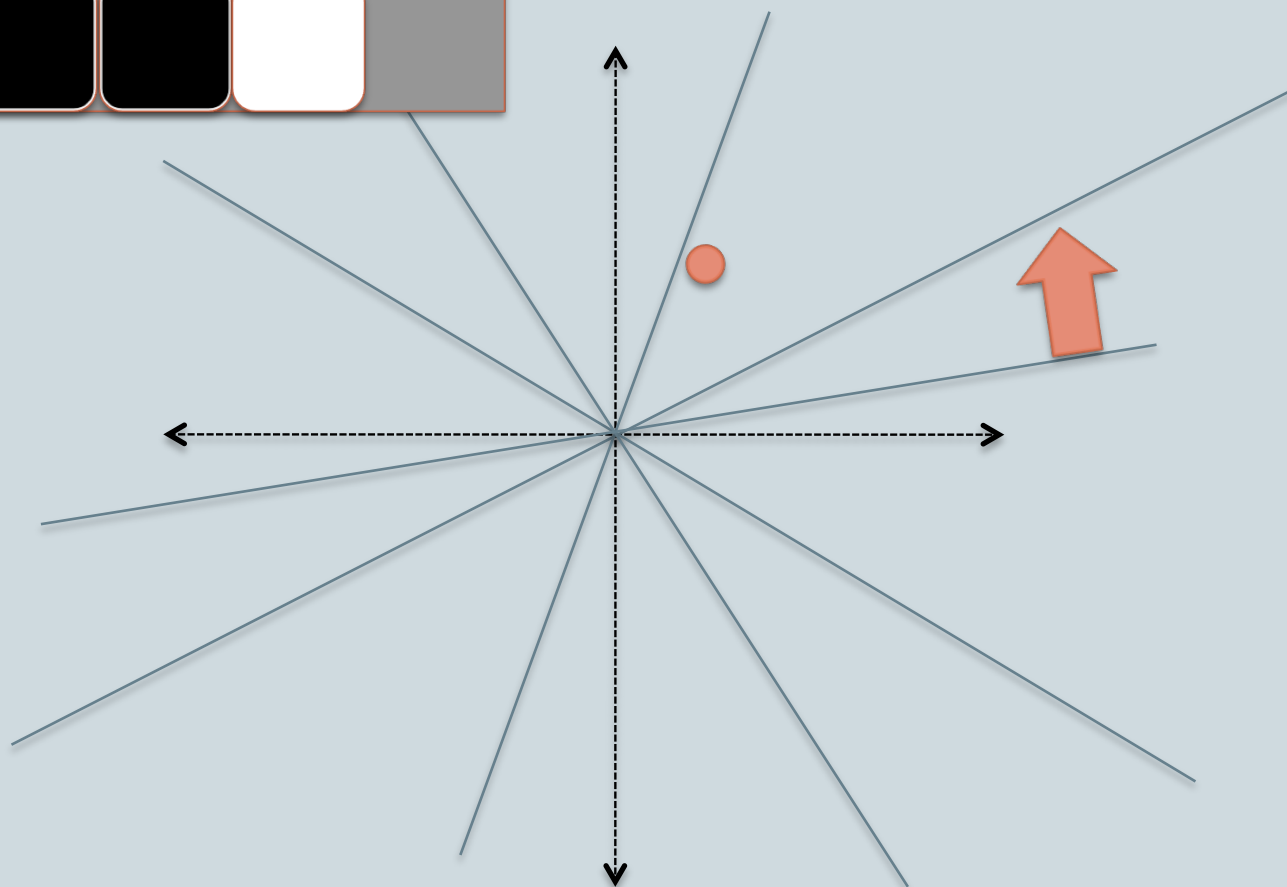
Locality Sensitive Hashing



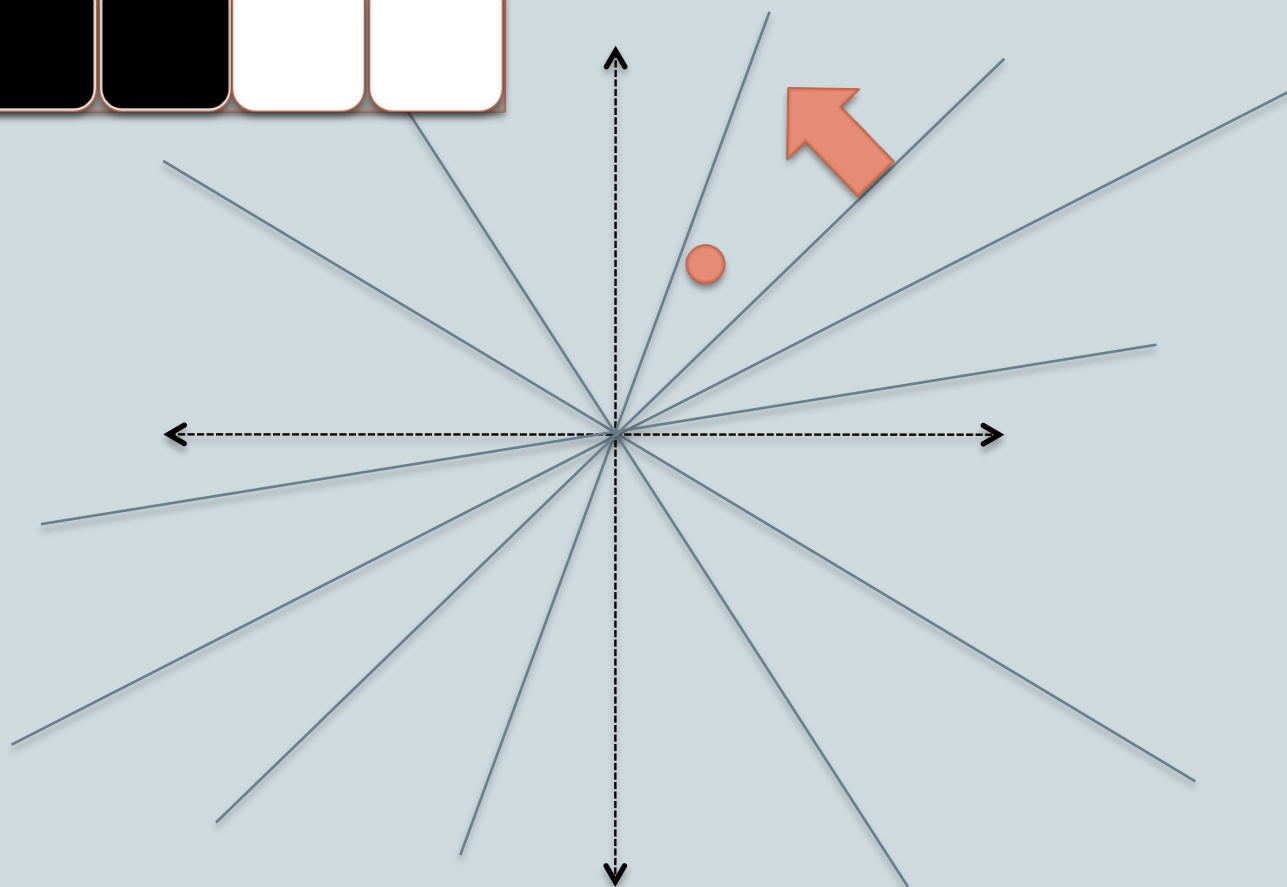
Locality Sensitive Hashing



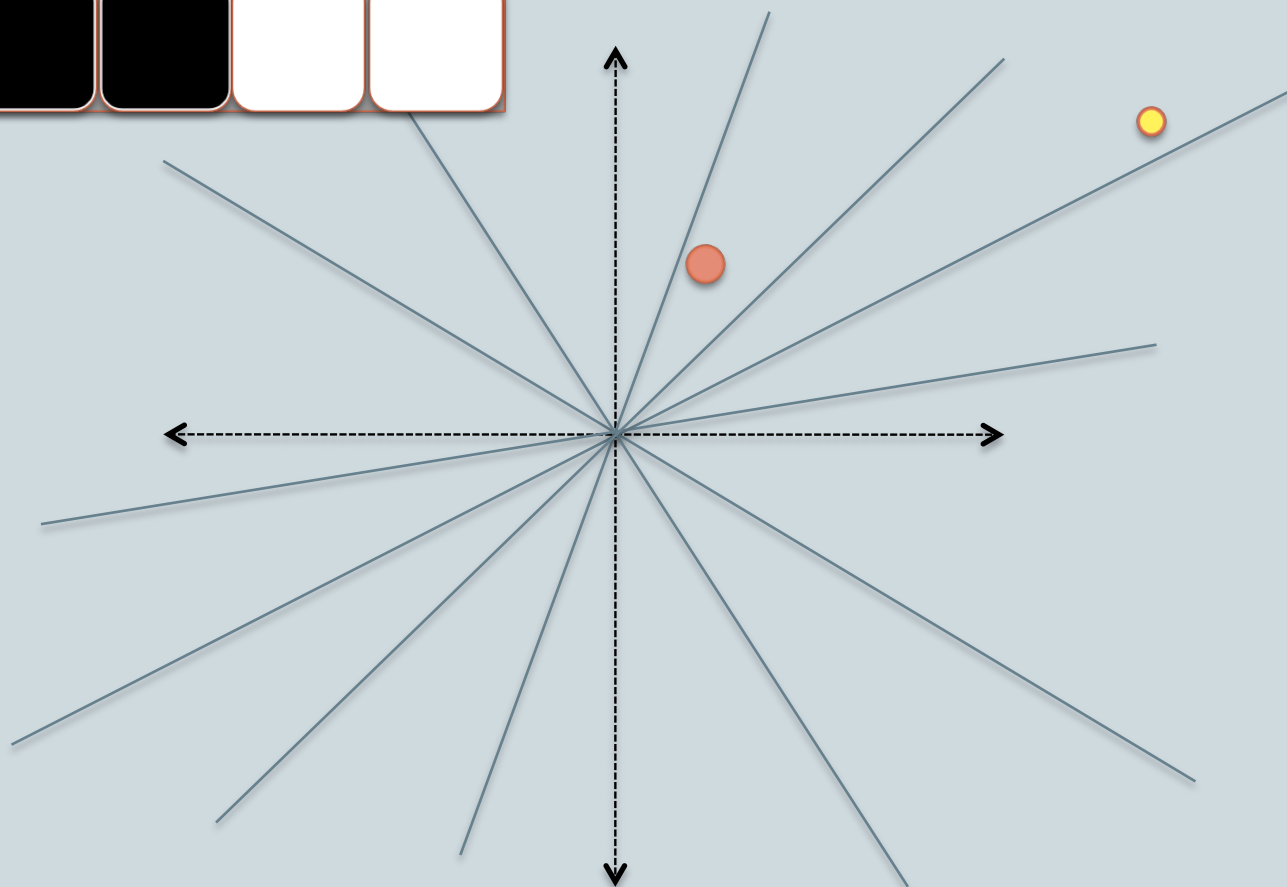
Locality Sensitive Hashing



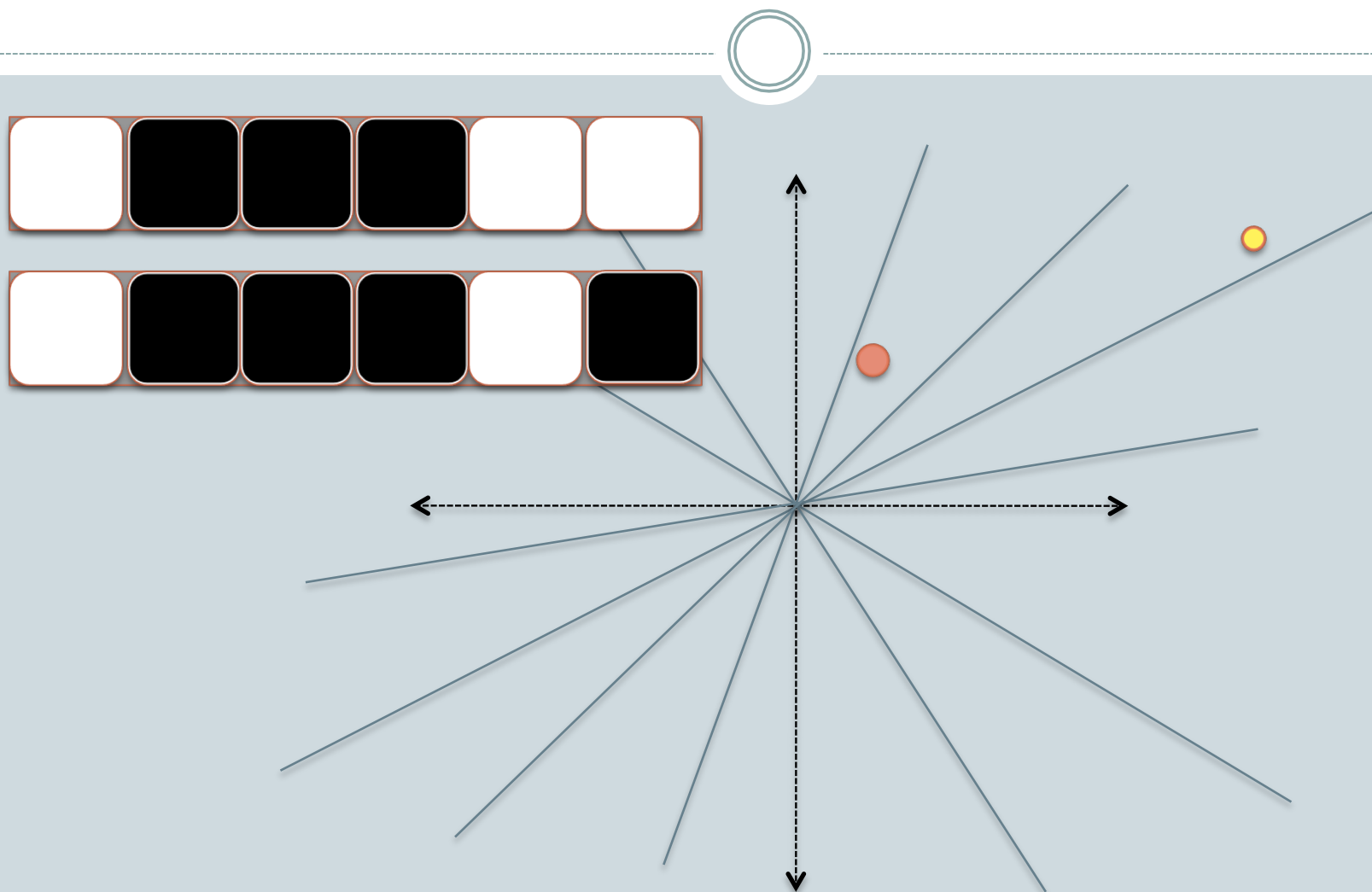
Locality Sensitive Hashing



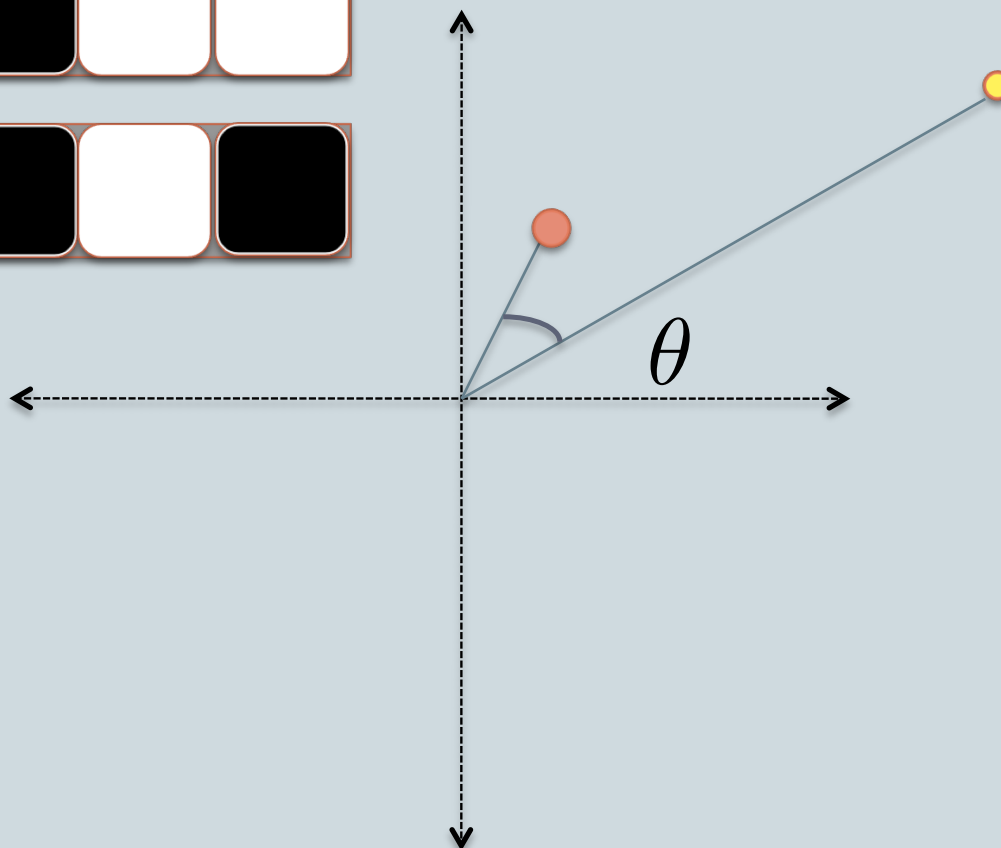
Locality Sensitive Hashing



Locality Sensitive Hashing



Locality Sensitive Hashing

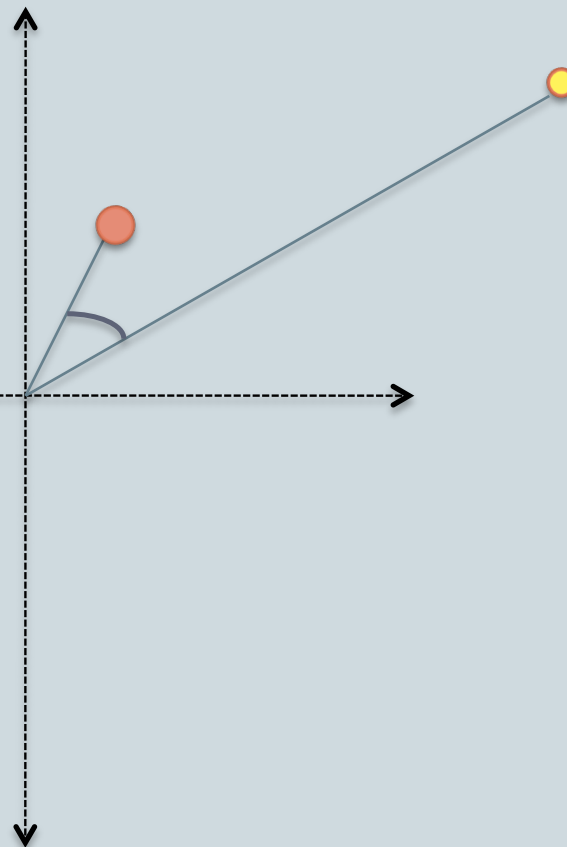


Locality Sensitive Hashing



Hamming Distance $\coloneqq h = 1$

Signature Length $\coloneqq b = 6$

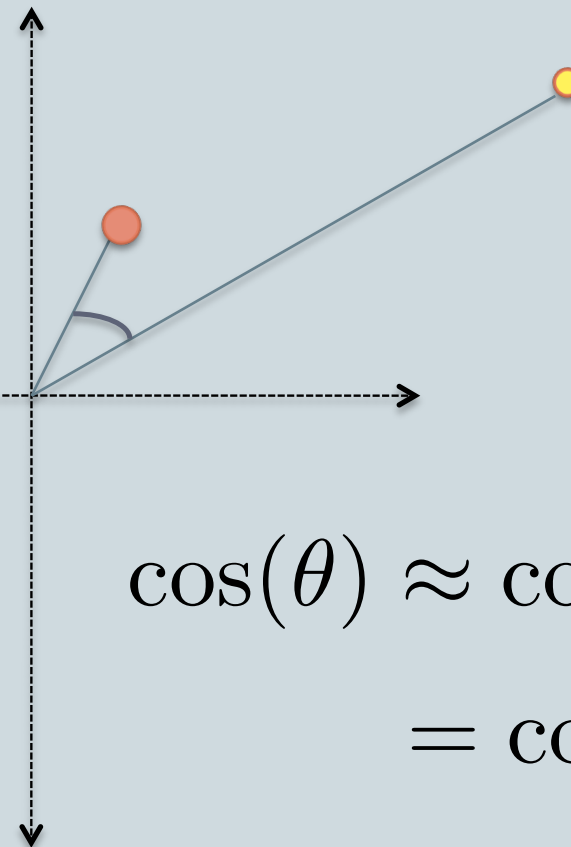


Locality Sensitive Hashing



Hamming Distance $:= h = 1$

Signature Length $:= b = 6$

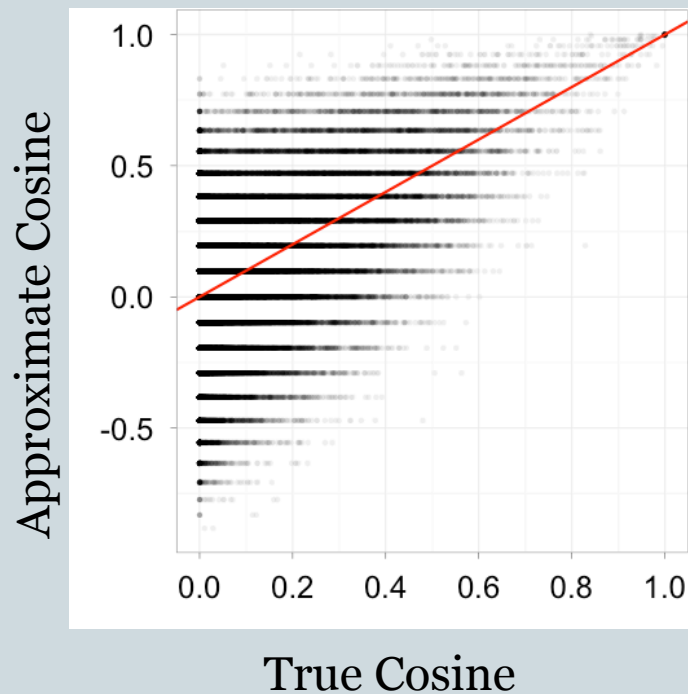


$$\begin{aligned}\cos(\theta) &\approx \cos\left(\frac{h}{b}\pi\right) \\ &= \cos\left(\frac{1}{6}\pi\right)\end{aligned}$$

Accuracy as function of bit length

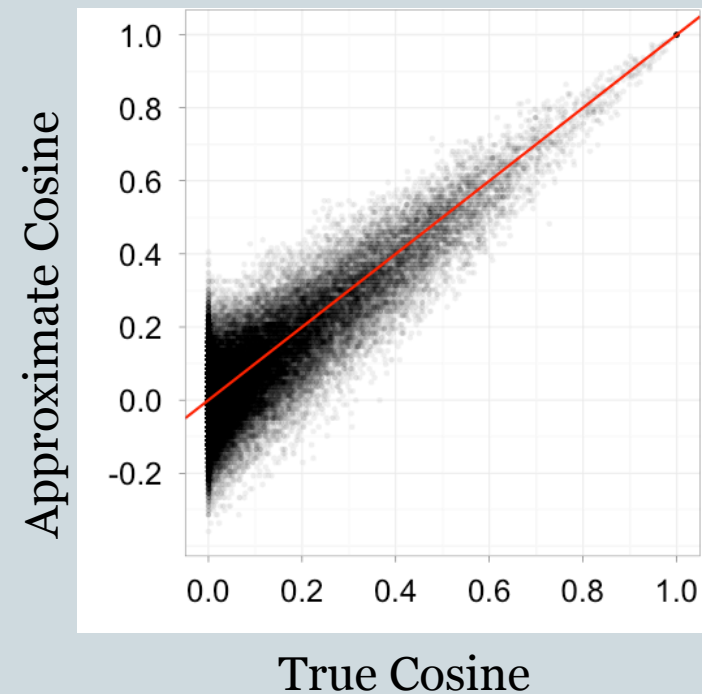


32 bit signatures



Cheap

256 bit signatures



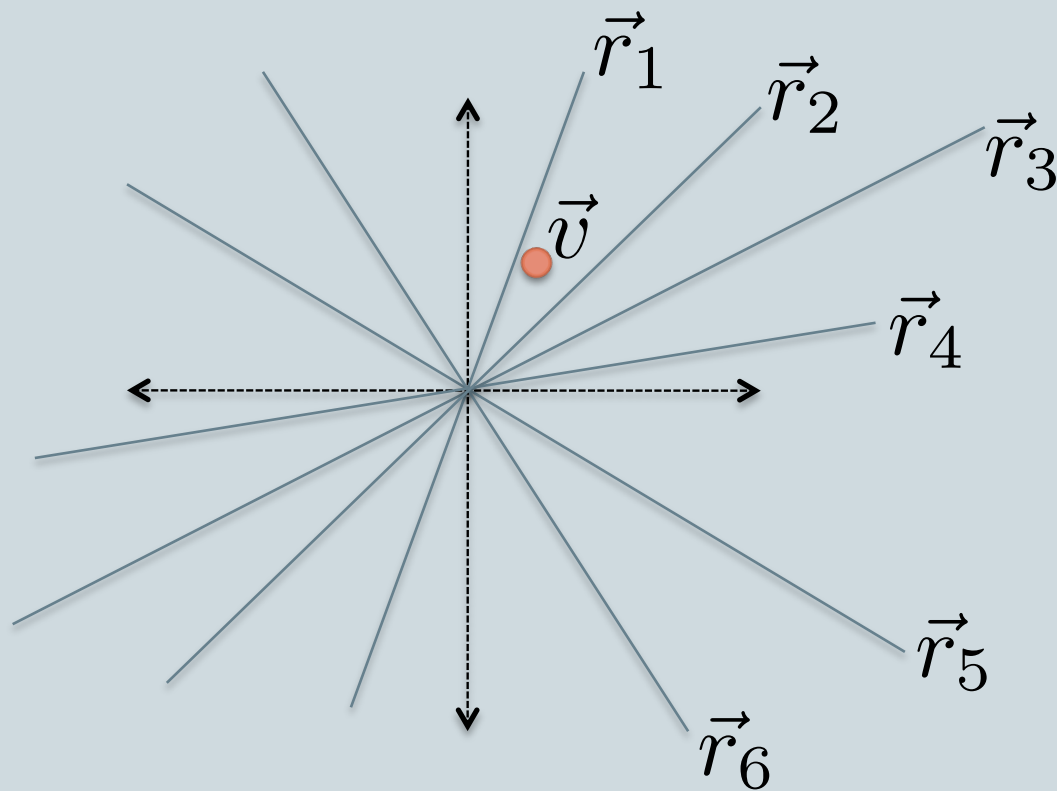
Accurate

Novel Points Here



1. Online hash function
2. Pooling trick

Online Hash Function



$$\vec{v} \in \mathbb{R}^d$$

$$\vec{r}_i \sim N(0, 1)^d$$

Online Hash Function



$$\vec{v} \in \mathbb{R}^d$$

$$\vec{r}_i \sim N(0, 1)^d$$

$$h_i(\vec{v}) = \begin{cases} 1 & \text{if } \vec{v} \cdot \vec{r}_i \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Online Hash Function



$$\begin{aligned} \text{if } \vec{v} &= \sum_j \vec{v}_j \\ \text{then } \vec{v} \cdot \vec{r}_i &= \sum_j \vec{v}_j \cdot \vec{r}_i \end{aligned}$$

Break into local products

Online Hash Function



$$h_i(\vec{v}) = \begin{cases} 1 & \text{if } \vec{v} \cdot \vec{r}_i \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Offline

if $\vec{v} = \sum_j \vec{v}_j$

then $\vec{v} \cdot \vec{r}_i = \sum_j \vec{v}_j \cdot \vec{r}_i$

Online

$$h_{it}(\vec{v}) = \begin{cases} 1 & \text{if } \sum_j^t \vec{v}_j \cdot \vec{r}_i \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

The Pooling Trick



$$\vec{r}_i \sim N(0, 1)^d$$

$$\begin{pmatrix} \vec{r}_1 \\ \dots \\ \vec{r}_b \end{pmatrix} = \begin{pmatrix} N(0, 1) & \dots & N(0, 1) \\ \dots & \dots & \dots \\ N(0, 1) & \dots & N(0, 1) \end{pmatrix}$$

The Pooling Trick



The random projection matrix can easily
require gigabytes of memory

$$\begin{pmatrix} \vec{r}_1 \\ \dots \\ \vec{r}_b \end{pmatrix} = \begin{pmatrix} N(0, 1) & \dots & N(0, 1) \\ \dots & \dots & \dots \\ N(0, 1) & \dots & N(0, 1) \end{pmatrix}$$

The Pooling Trick



$$\vec{p} \sim N(0, 1)^m$$

$$\vec{p} = (N(0, 1) \quad \dots \quad N(0, 1))$$

Define $\vec{r}_i[j] = \vec{p}[\text{hash}(i, j) \text{ MOD } m]$

$$m \ll b \times d$$

The Pooling Trick

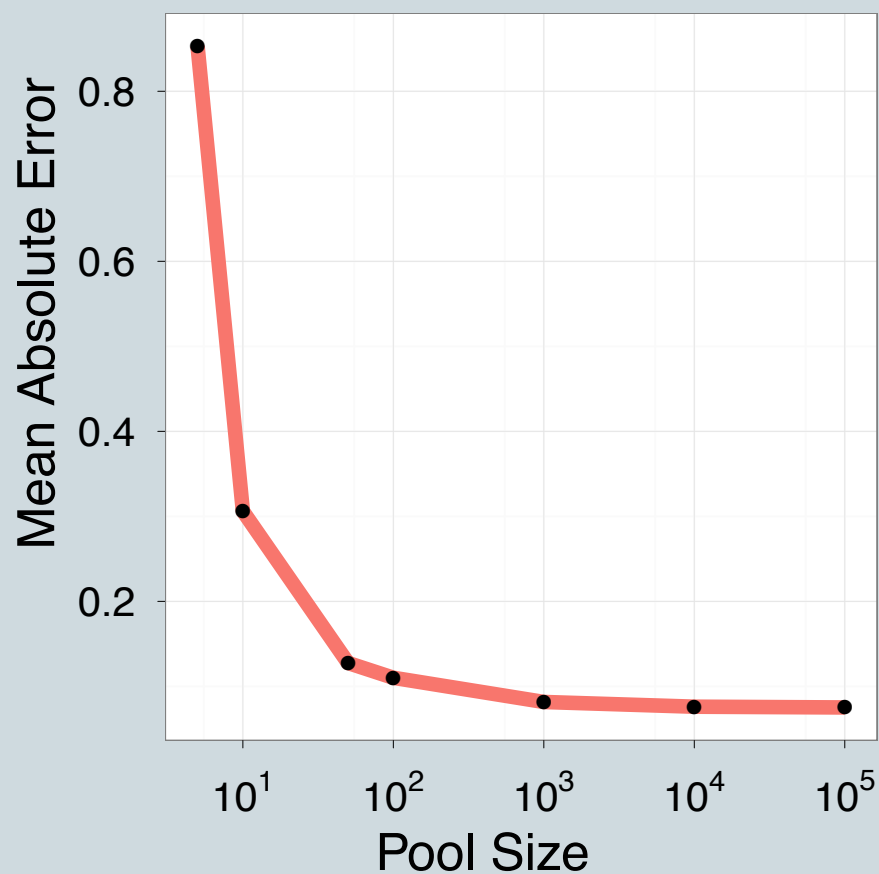


Define $\vec{r}_i[j] = \vec{p}[\text{hash}(i, j) \bmod m]$

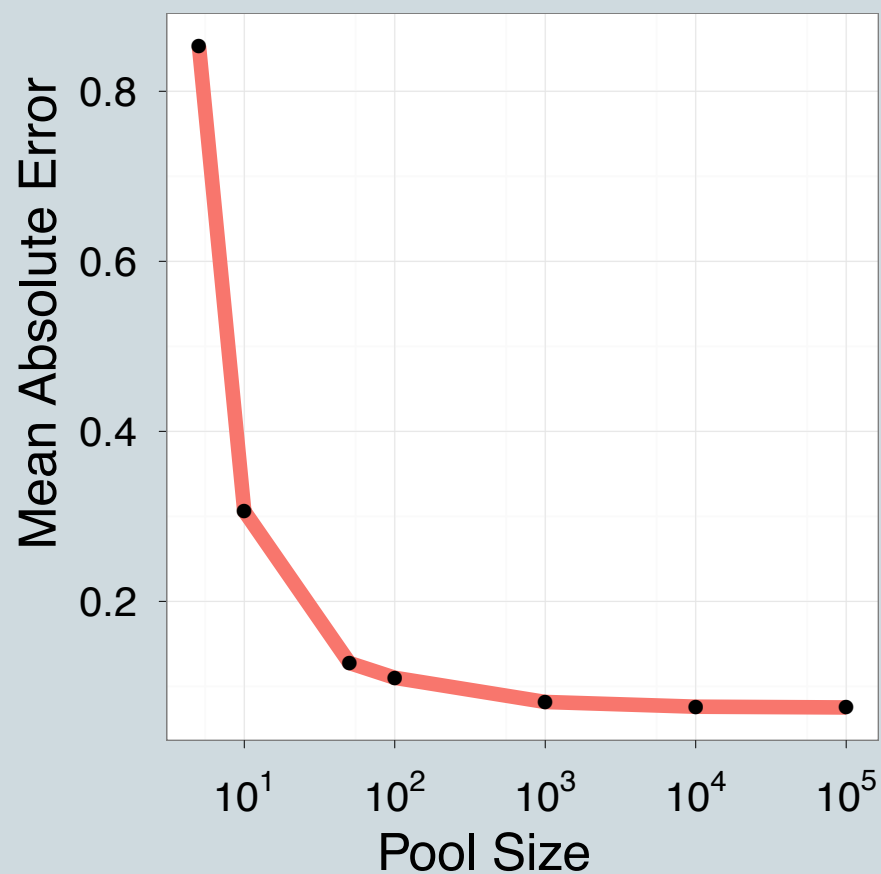
$$\begin{pmatrix} \dots & \dots & \dots \\ \dots & (i, j) & \dots \\ \dots & \dots & \dots \end{pmatrix}$$

$$\left(N(0, 1) \quad \dots \quad N(0, 1) \right)$$

The Pooling Trick



The Pooling Trick



That is 400 **kilobytes**

Example



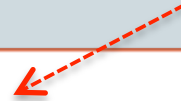
London

Milan.₉₇, **Madrid**.₉₆, **Stockholm**.₉₆, **Manila**.₉₅, **Moscow**.₉₅
ASHER₀, Champaign₀, MANS₀, NOBLE₀, come₀
Prague₁, Vienna₁, suburban₁, synchronism₁, Copenhagen₂
Frankfurt₄, Prague₄, Taszar₅, Brussels₆, Copenhagen₆
Prague₁₂, Stockholm₁₂, Frankfurt₁₄, Madrid₁₄, Manila₁₄
Stockholm₂₀, Milan₂₂, Madrid₂₄, Taipei₂₄, Frankfurt₂₅

Accurate



Closest based on true cosine



London

Milan_{.97}, Madrid_{.96}, Stockholm_{.96}, Manila_{.95}, Moscow_{.95}

ASHER₀, Champaign₀, MANS₀, NOBLE₀, come₀

Prague₁, Vienna₁, suburban₁, synchronism₁, Copenhagen₂

Frankfurt₄, Prague₄, Taszar₅, Brussels₆, Copenhagen₆

Prague₁₂, Stockholm₁₂, Frankfurt₁₄, Madrid₁₄, Manila₁₄

Stockholm₂₀, Milan₂₂, Madrid₂₄, Taipei₂₄, Frankfurt₂₅



London

Milan.₉₇, Madrid.₉₆, Stockholm.₉₆, Manila.₉₅, Moscow.₉₅

ASHER₀, Champaign₀, MANS₀, NOBLE₀, come₀

Prague₁, Vienna₁, suburban₁, synchronism₁, Copenhagen₂

Frankfurt₄, Prague₄, Taszar₅, Brussels₆, Copenhagen₆

Prague₁₂, Stockholm₁₂, Frankfurt₁₄, Madrid₁₄, Manila₁₄

Stockholm₂₀, Milan₂₂, Madrid₂₄, Taipei₂₄, Frankfurt₂₅

Closest based on 32 bit sig.'s

Cheap



London

Milan.₉₇, Madrid.₉₆, Stockholm.₉₆, Manila.₉₅, Moscow.₉₅
ASHER₀, Champaign₀, MANS₀, NOBLE₀, come₀
Prague₁, Vienna₁, suburban₁, synchronism₁, Copenhagen₂
Frankfurt₄, Prague₄, Taszar₅, Brussels₆, Copenhagen₆
~~Prague₁₂, Stockholm₁₂, Frankfurt₁₄, Madrid₁₄, Manila₁₄~~
Stockholm₂₀, Milan₂₂, Madrid₂₄, Taipei₂₄, Frankfurt₂₅

Closest based on 256 bit sig.'s

Cheap-ish

No Questions.



See Poster.

Also: these algorithms distribute to the cloud.



JOHNS HOPKINS
U N I V E R S I T Y

