# Hash Kernels & Linear Models

# Vector Space Model for Text

## A Vector Space Model for Automatic Indexing

G. Salton, A. Wong
and C. S. Yang
Cornell University

In a document retrieval, or other pattern matching environment where stored entities (documents) are compared with each other or with incoming patterns (search requests), it appears that the best indexing (property) space is one where each entity lies as far away from the others as possible; in these circumstances the value of an indexing system may be expressible as a function of the density of the object space; in particular, retrieval performance may correlate inversely with space density. An approach based on space density computations is used to choose an optimum indexing vocabulary for a collection of documents. Typical evaluation results are shown, demonstrating the usefulness of the model.
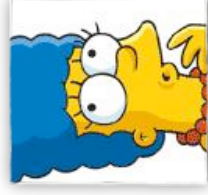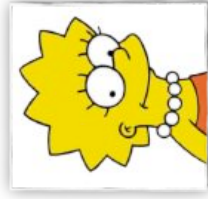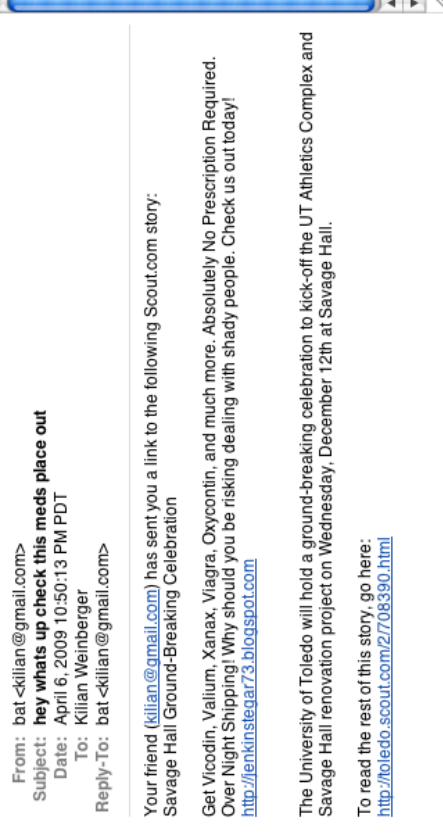
Carnegie Mellon University

# Linear functions

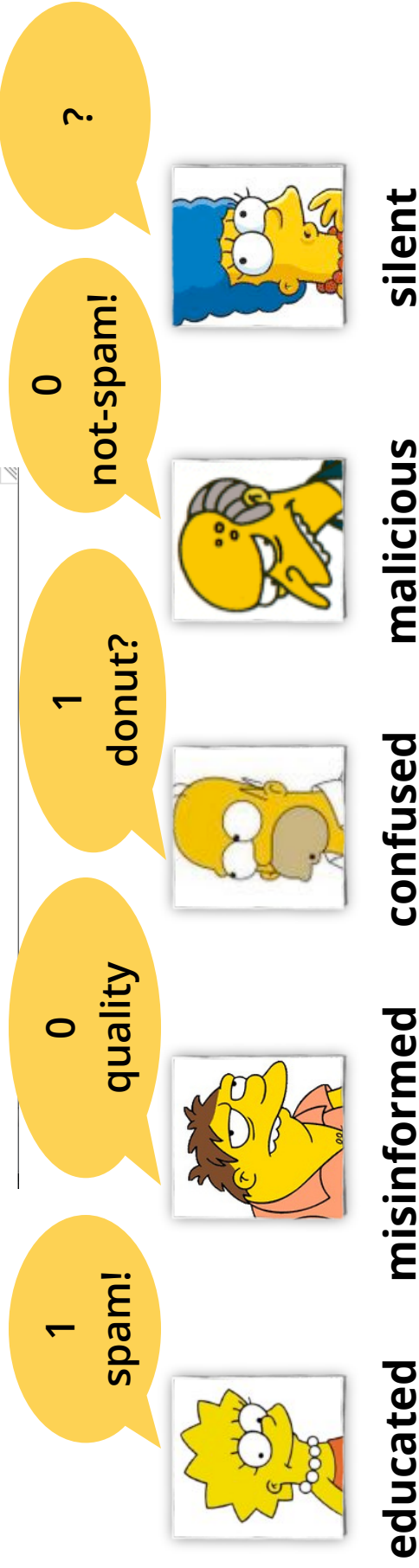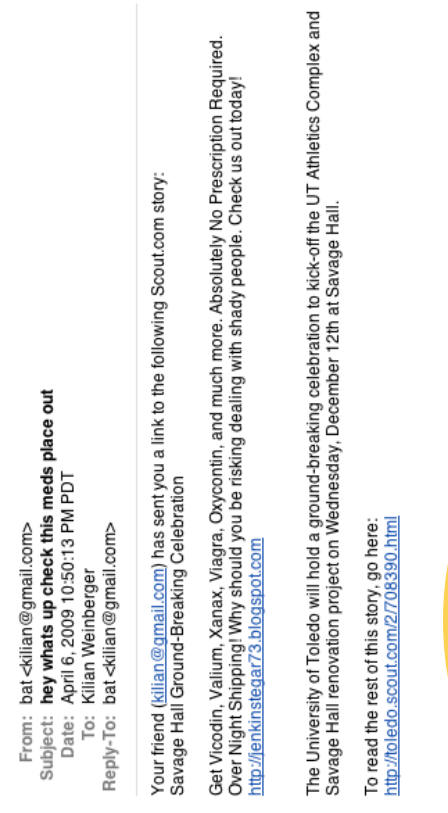- Function

$$f(x) = \langle w, x \rangle + b = \sum_i w_i x_i + b$$

- Used for spam filtering, internet advertising, ranking, retrieval, summarization, gene finding, stock prediction, user profiling, ...

- Good as long as we don't have more than 1 billion parameters (8GB of memory for double precision)

- What if we have more parameters?

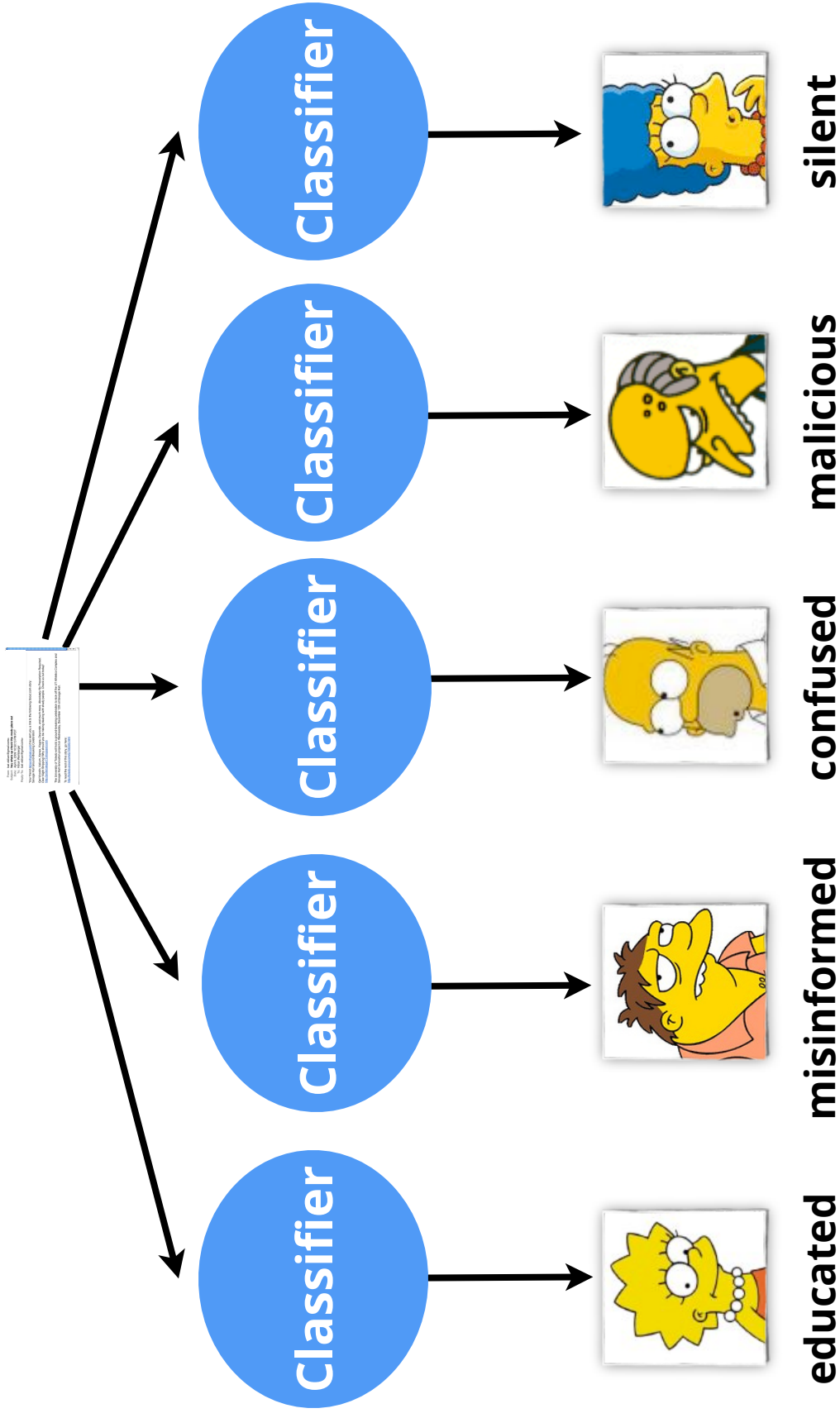  - Lasso (with distributed optimization)

  - Hashing

# Personalized Spam Filtering

From: bat <kilian@gmail.com>
Subject: **hey whats up check this meds place out**
Date: April 6, 2009 10:50:13 PM PDT
To: Kilian Weinberger
Reply-To: bat <kilian@gmail.com>

Your friend (kilian@gmail.com) has sent you a link to the following Scout.com story:
Savage Hall Ground-Breaking Celebration

Get Vicodin, Valium, Xanax, Viagra, Oxycontin, and much more. Absolutely No Prescription Required.
Over Night Shipping! Why should you be risking dealing with shady people. Check us out today!
http://fenkinstegar73.blogspot.com

The University of Toledo will hold a ground-breaking celebration to kick-off the UT Athletics Complex and
Savage Hall renovation project on Wednesday, December 12th at Savage Hall.

To read the rest of this story, go here:
http://toledo.scout.com/2/708390.html

# Personalized Spam Filtering

From: bat <kilian@gmail.com>
Subject: **hey whats up check this meds place out**
Date: April 6, 2009 10:50:13 PM PDT
To: Kilian Weinberger
Reply-To: bat <kilian@gmail.com>

Your friend (kilian@gmail.com) has sent you a link to the following Scout.com story:
Savage Hall Ground-Breaking Celebration

Get Vicodin, Valium, Xanax, Viagra, Oxycontin, and much more. Absolutely No Prescription Required.
Over Night Shipping! Why should you be risking dealing with shady people. Check us out today!
http://jenkinstegar73.blogspot.com

The University of Toledo will hold a ground-breaking celebration to kick-off the UT Athletics Complex and
Savage Hall renovation project on Wednesday, December 12th at Savage Hall.

To read the rest of this story, go here:
http://toledo.scout.com/2/708390.html

**1 spam!** **0 quality** **1 donut?** **0 not-spam!** **?**

**educated** **misinformed** **confused** **malicious** **silent**

# Personalized Spam Filtering

# Multitask Learning

**Global Classifier**

Classifier — silent

Classifier — malicious

Classifier — confused

Classifier — misinformed

Classifier — educated

# Collaborative Classification

- **Primal space representation**

$$f(x, u) = \langle \phi(x), w \rangle + \langle \phi(x), w_u \rangle = \langle \phi(x) \otimes (1 \oplus e_u), w \rangle$$
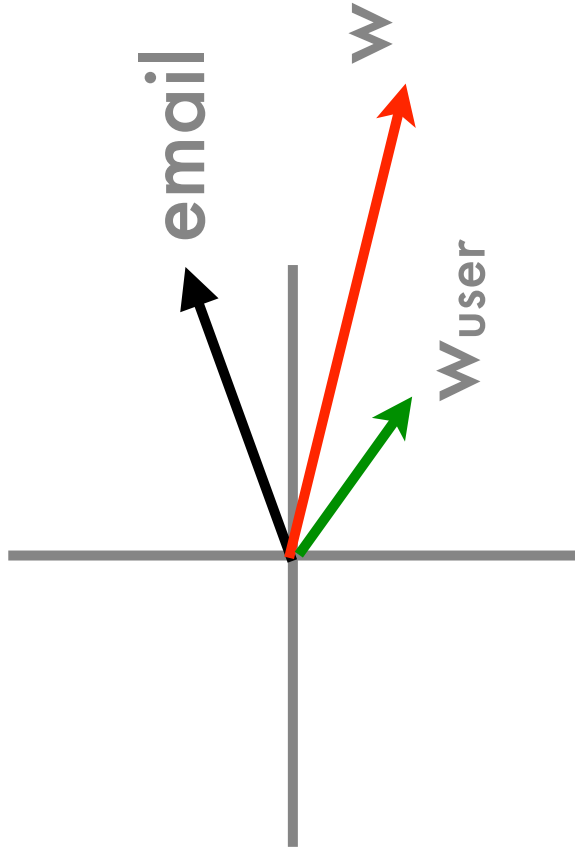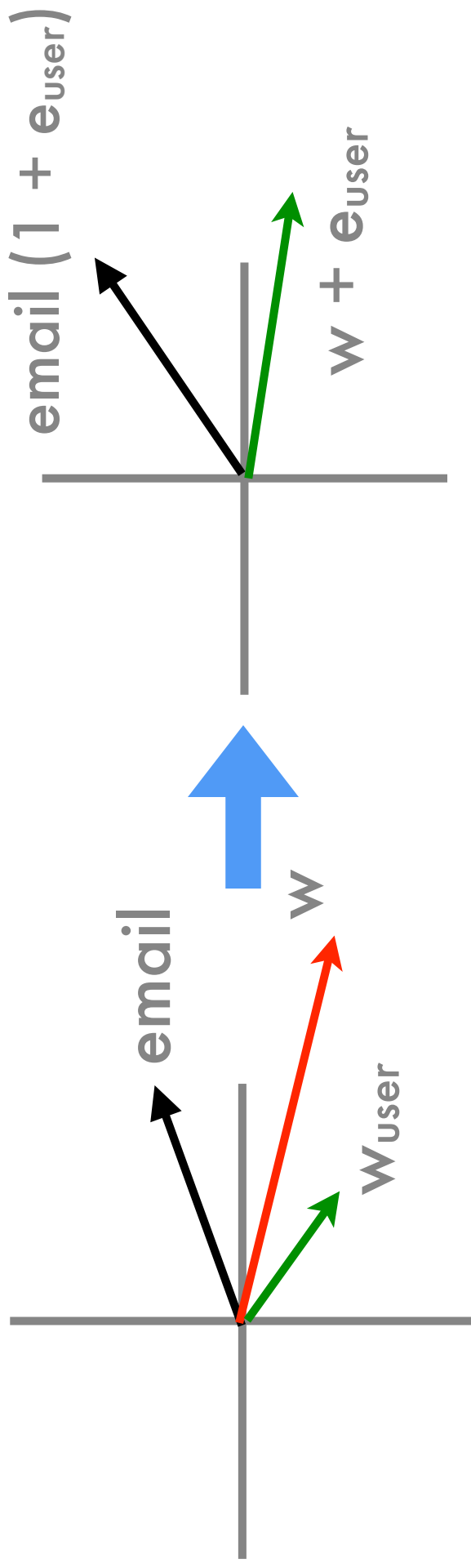
  **Kernel representation**

$$k((x, u), (x', u')) = k(x, x')[1 + \delta_{u, u'}]$$

  Multitask kernel (e.g. Pontil & Michelli, Daume). Usually does not scale well ...

- **Problem** - dimensionality is $10^{13}$. That is 40TB of space

# Collaborative Classification



- **Primal space representation**

$$f(x, u) = \langle \phi(x), w \rangle + \langle \phi(x), w_u \rangle = \langle \phi(x) \otimes (1 \oplus e_u), w \rangle$$
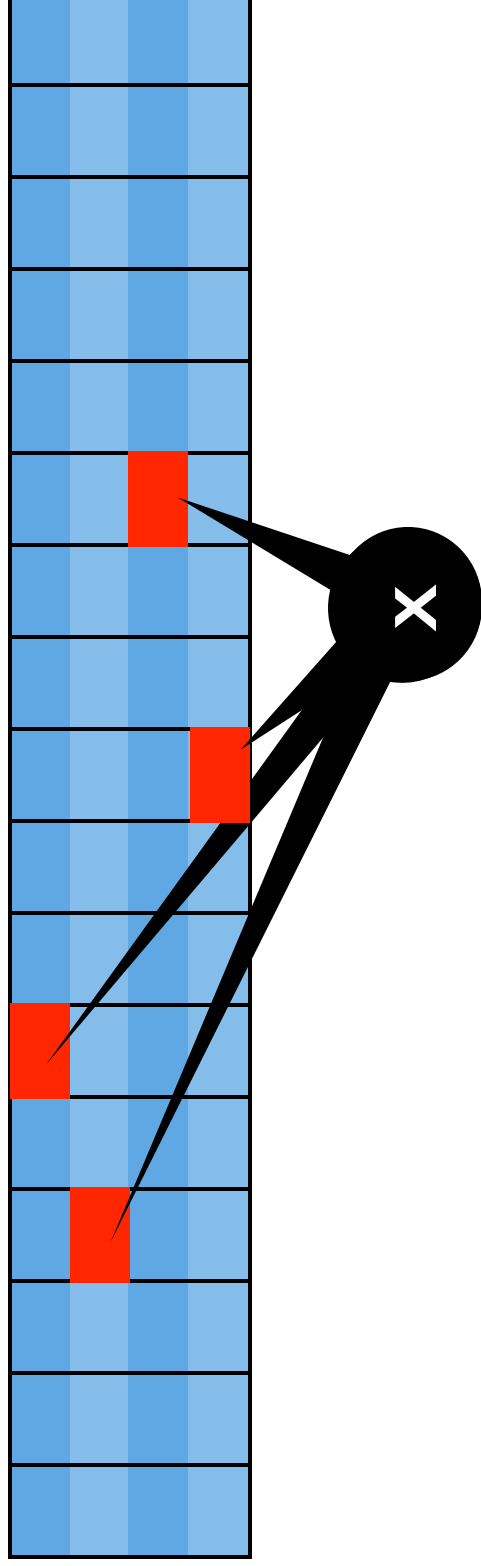
  **Kernel representation**

$$k((x, u), (x', u')) = k(x, x')[1 + \delta_{u, u'}]$$

  Multitask kernel (e.g. Pontil & Michelli, Daume). Usually does not scale well ...
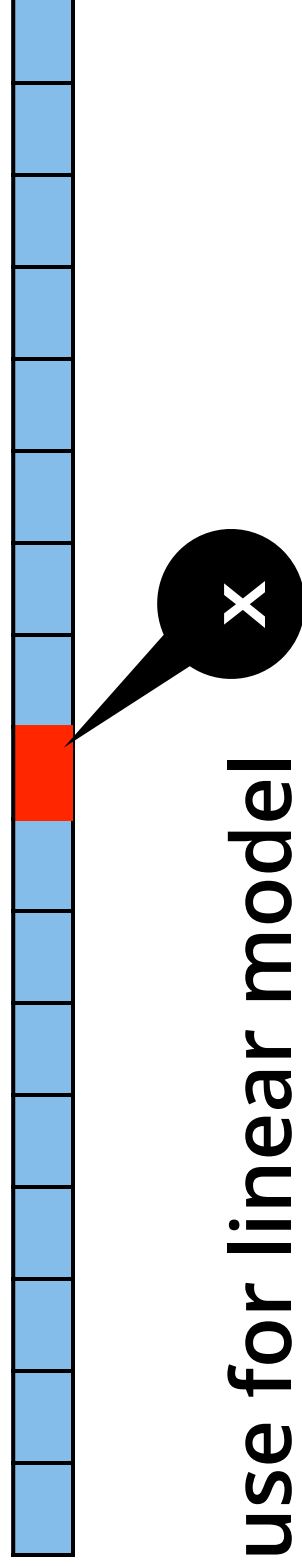
- **Problem** - dimensionality is $10^{13}$. That is 40TB of space

# Collaborative Classification



- **Primal space representation**

$$f(x, u) = \langle \phi(x), w \rangle + \langle \phi(x), w_u \rangle = \langle \phi(x) \otimes (1 \oplus e_u), w \rangle$$

- **Kernel representation**

$$k((x, u), (x', u')) = k(x, x')[1 + \delta_{u,u'}]$$

Multitask kernel (e.g. Pontil & Michelli, Daume). Usually does not scale well ...

- **Problem** - dimensionality is $10^{13}$. That is 40TB of space

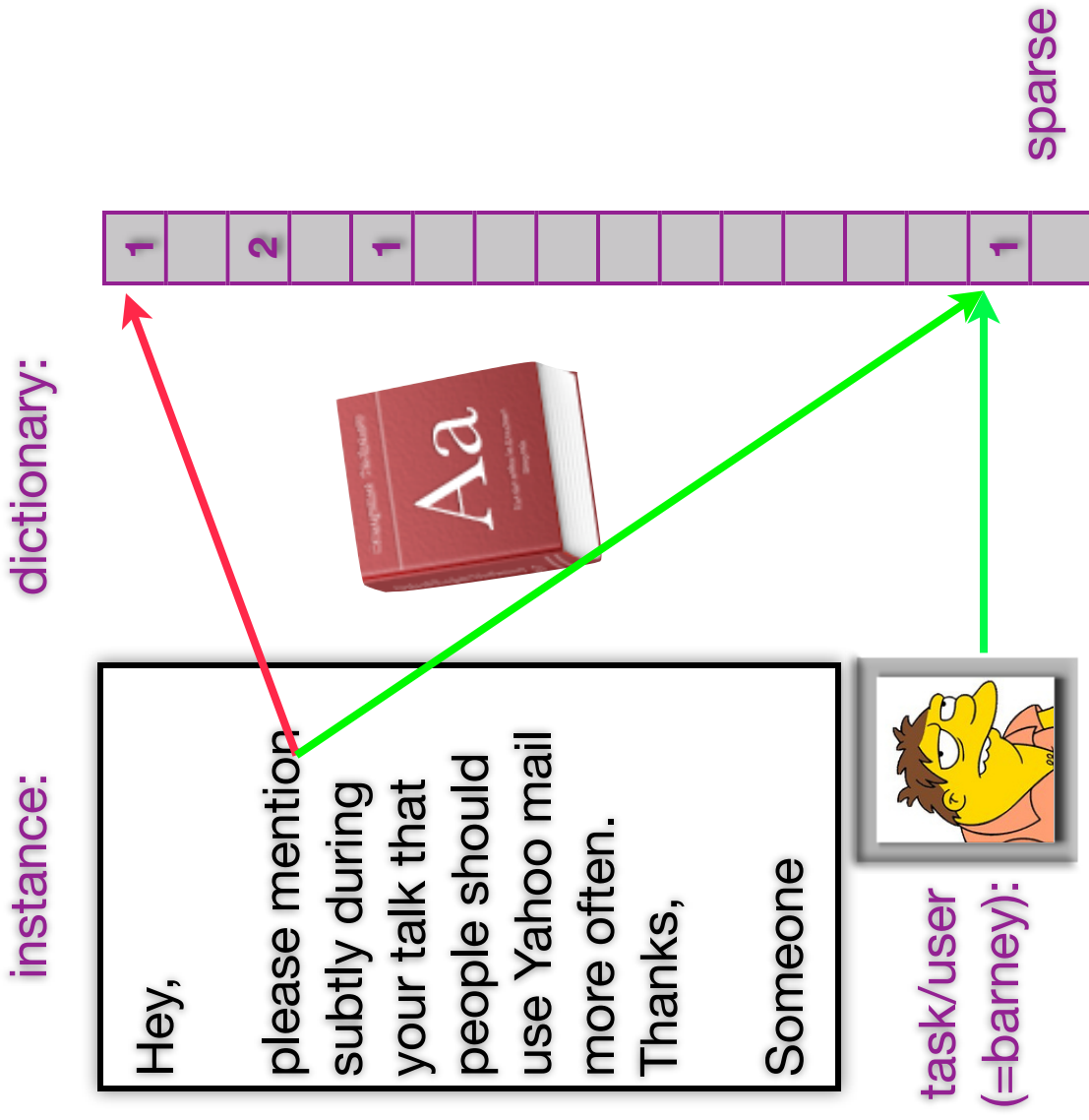# CountMin Datastructure

- Count Sketch Basic
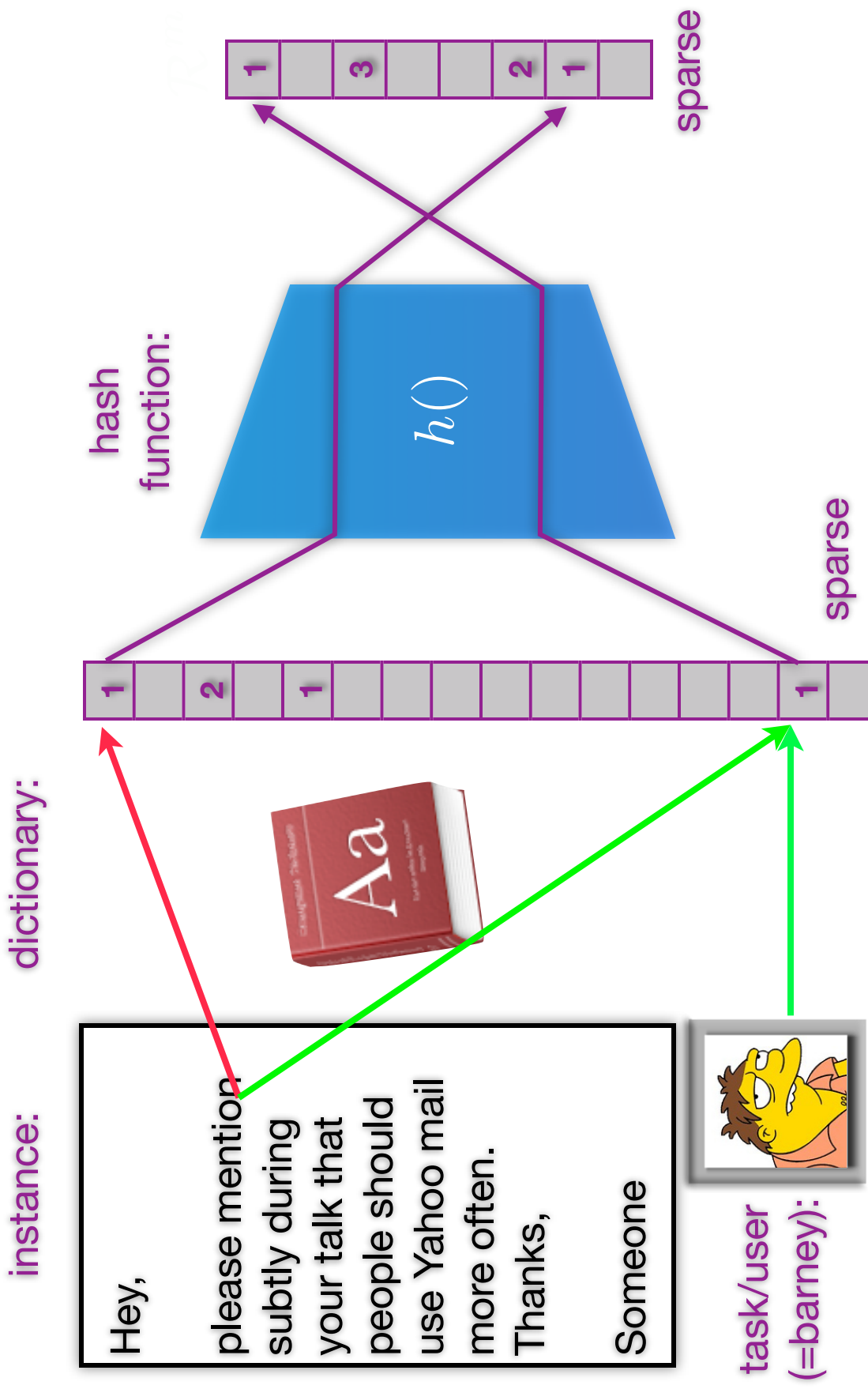
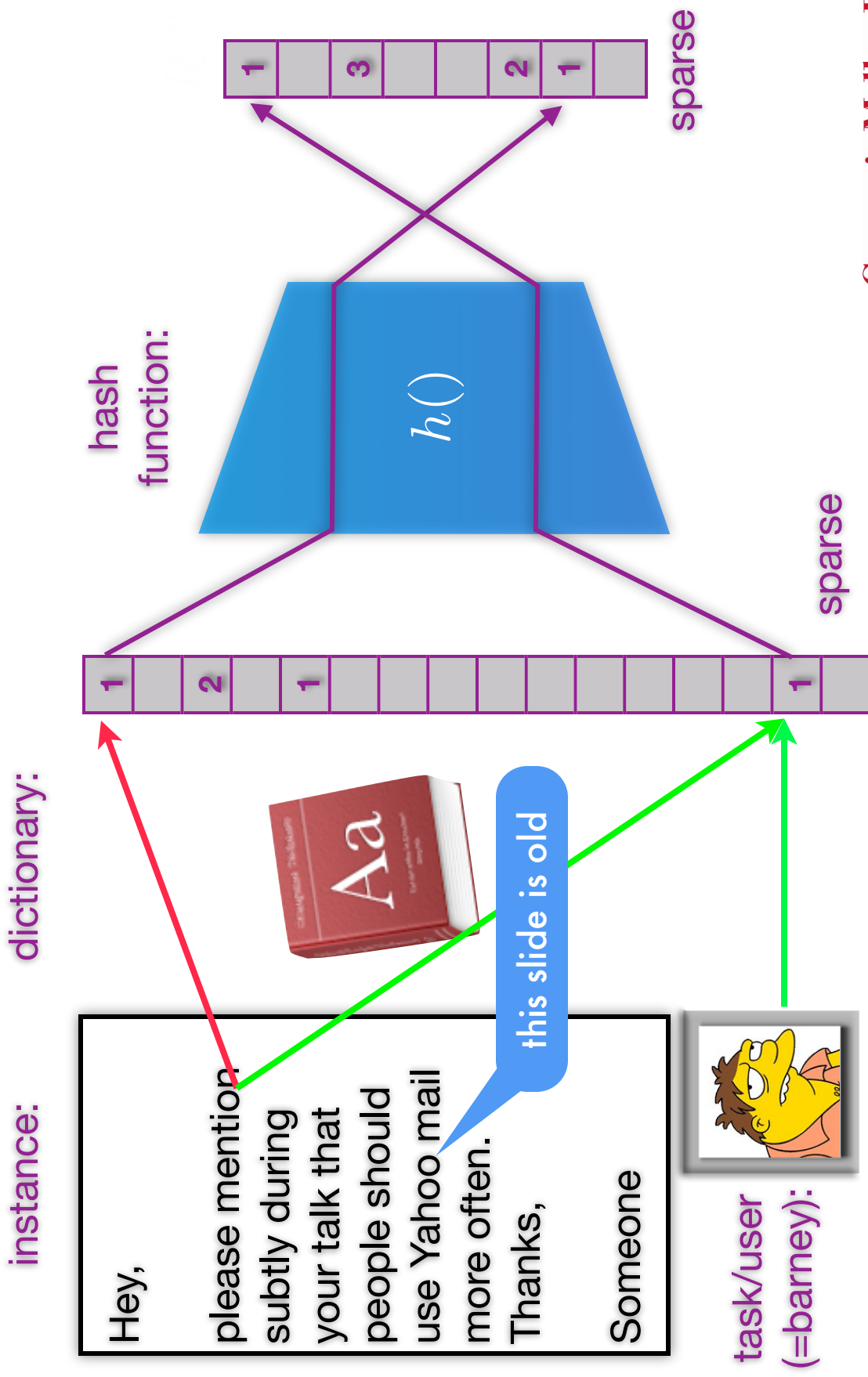- Making it even more sparse

use for linear model

# Hash Kernel

# Hash Kernel

dictionary:

instance:

Hey,

please mention
subtly during
your talk that
people should
use Yahoo mail
more often.
Thanks,

Someone

task/user
(=barney):

| 1 | | 2 | 1 | | | | | | | | | 1 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

sparse

# Hash Kernel

instance:

Hey,

please mention subtly during your talk that people should use Yahoo mail more often. Thanks,
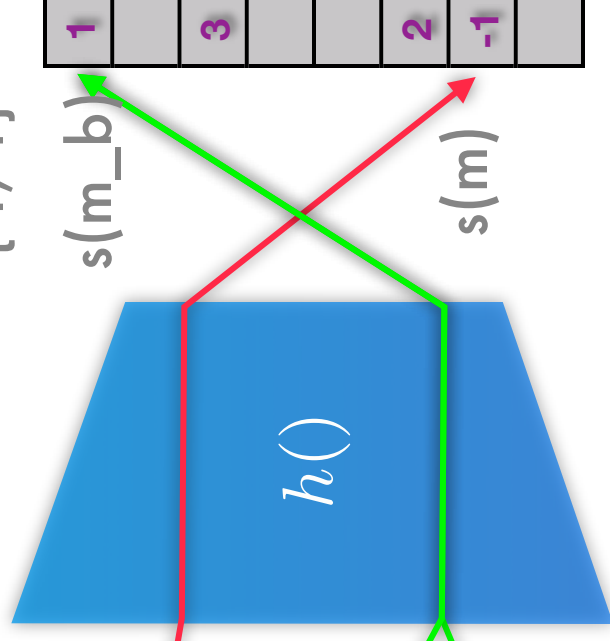
Someone

task/user (=barney):

dictionary:

hash function: $h()$

sparse

| 1 | | 2 | | 1 | | | | | | | | | | 1 | |

sparse

| 1 | | 3 | | | 2 | 1 | |

# Hash Kernel

$$f(x) = \sum_i w[h(i)]\sigma(i)x_i$$

{-1, 1}

instance:

Hey,

please mention
subtly during
your talk that
people should
use Yahoo mail
more often.
Thanks,

Someone

task/user
(=barney):

h('mention')

h('mention_barney')

$h()$

s(m_b)
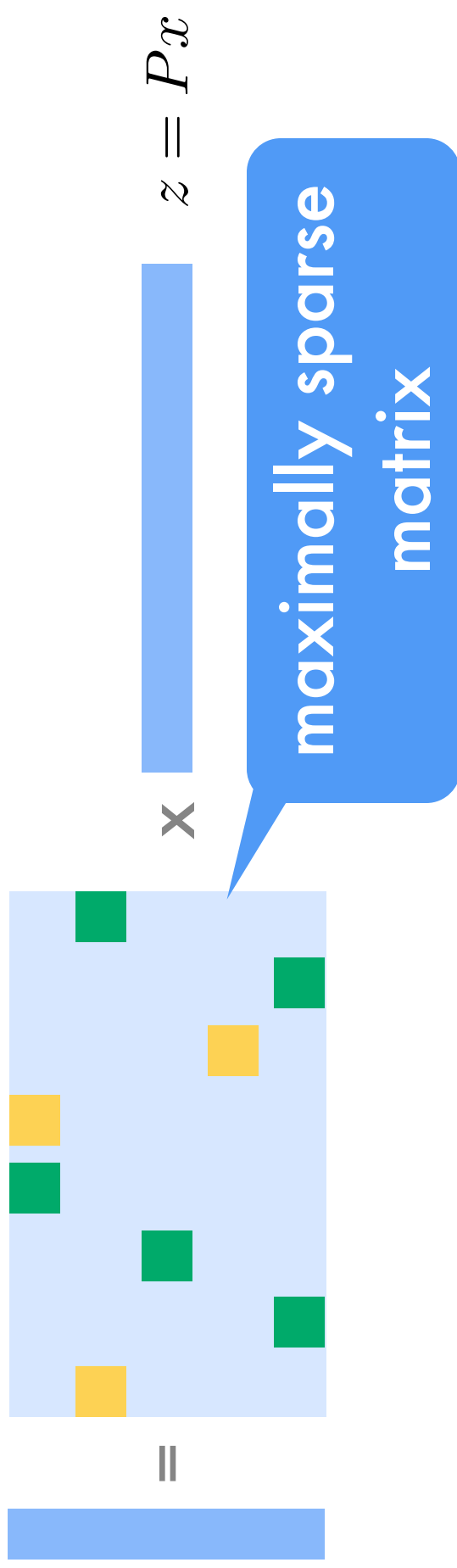
s(m)

| 1 | | 3 | | 2 | -1 | |
|---|---|---|---|---|---|---|

Similar to count hash

(Charikar, Chen, Farrach-Colton, 2003)

# Advantages of hashing

- No dictionary!
  - Content drift is no problem
  - All memory used for classification
  - Finite memory guarantee
    (good for online learning)
  - No Memory needed for projection (vs. LSH).
  - Implicit mapping into high dimensional space!
  - Sparsity preserving! (vs LSH)
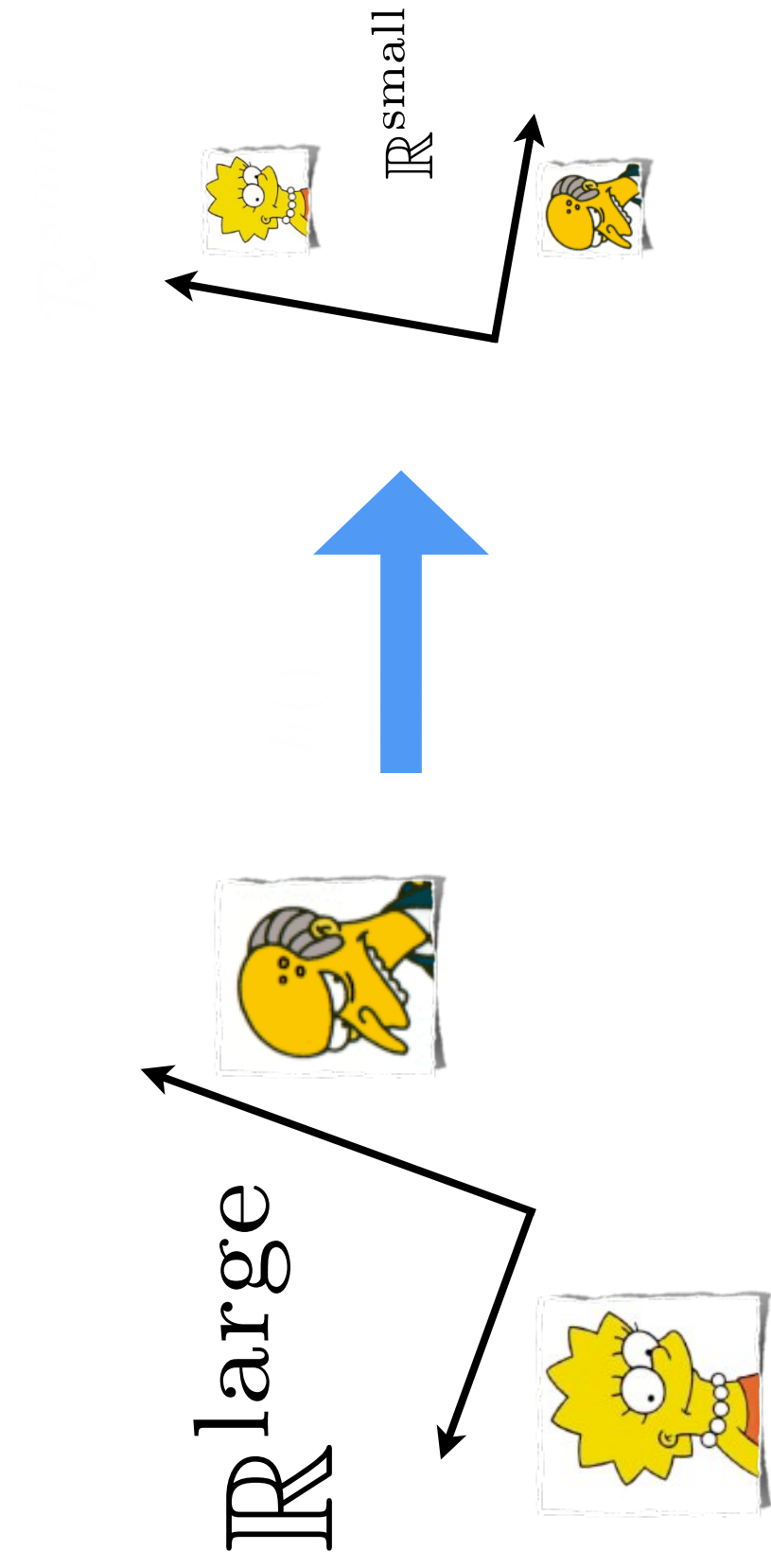
# Hash Kernels - the matrix view



$z = Px$

maximally sparse matrix

- Preserves inner product

$$\langle w, x \rangle = \sum_i w_i x_i \qquad \langle \bar{w}, \bar{x} \rangle = \sum_j \left[ \sum_{i:h(i)=j} w_i \sigma(i) \right] \left[ \sum_{i:h(i)=j} x_i \sigma(i) \right]$$

## Rademacher hash

$$\mathbb{E}_\sigma [\sigma(i)\sigma(i')] = \delta_{ii'}$$

# Approximate Orthogonality
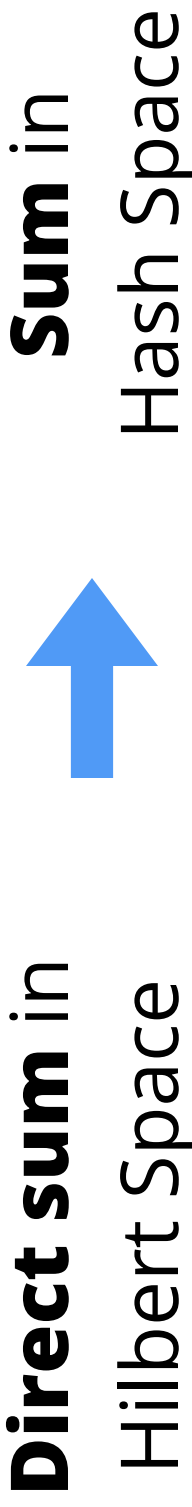
$\mathbb{R}^{large}$

$\mathbb{R}^{small}$

We can do multi-task learning!

# Guarantees

- For a random hash function the inner product vanishes with high probability via

$$\Pr\{|\langle w_v, h_u(x)\rangle)| > \epsilon\} \leq 2e^{-C\epsilon^2 m}$$
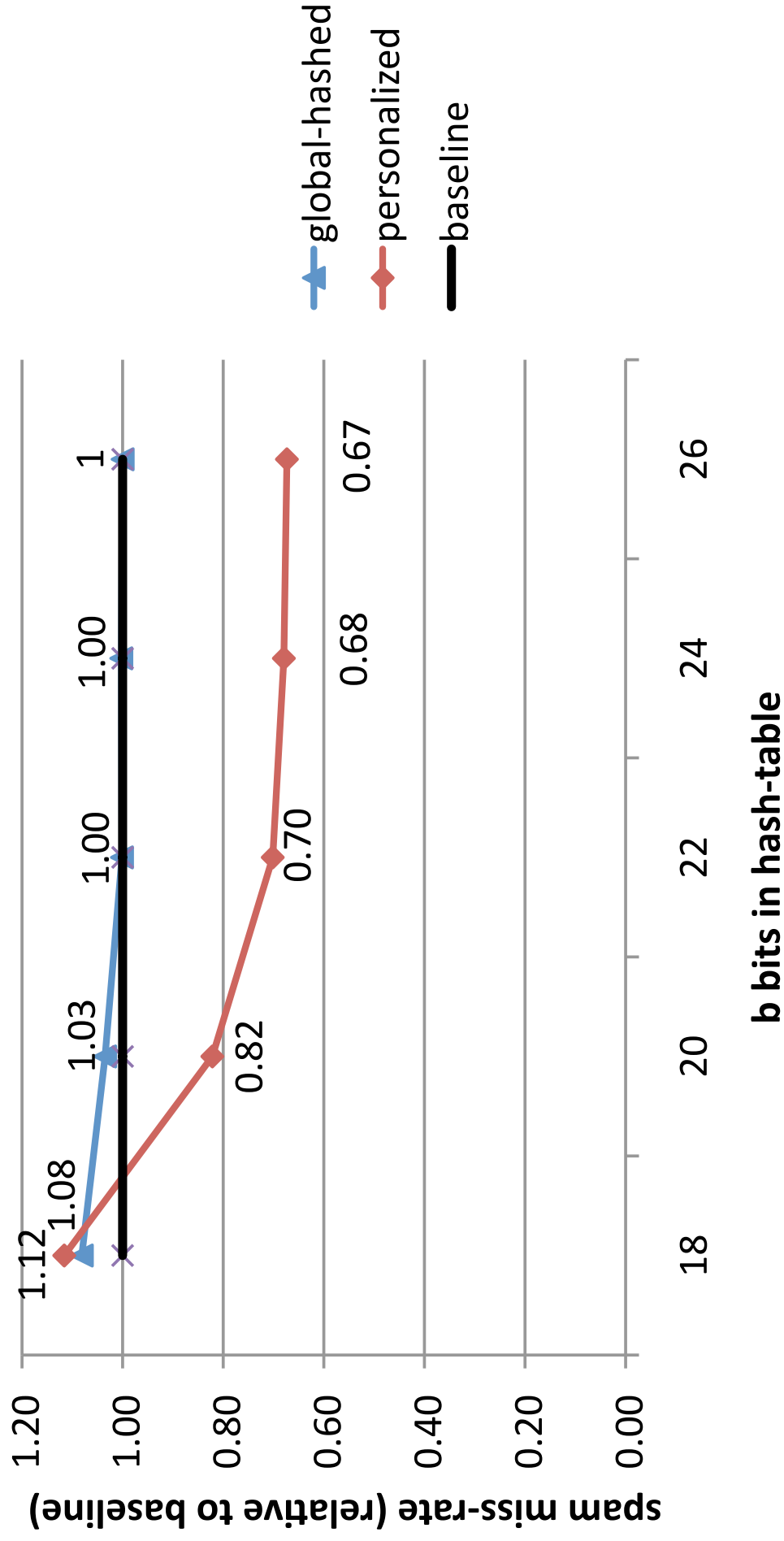
- We can use this for multitask learning

**Direct sum** in Hilbert Space



**Sum** in Hash Space

- Hashed inner product is unbiased
- Variance is O(1/n)
- Restricted isometry property (Kumar, Sarlos, Dasgupta 2010)

Weinberger, K., Dasgupta, A., Attenberg, J., Langford, J., and Smola, A. J., Feature Hashing for Large Scale Multitask Learning, International Conference on Machine Learning, 2009. PDF

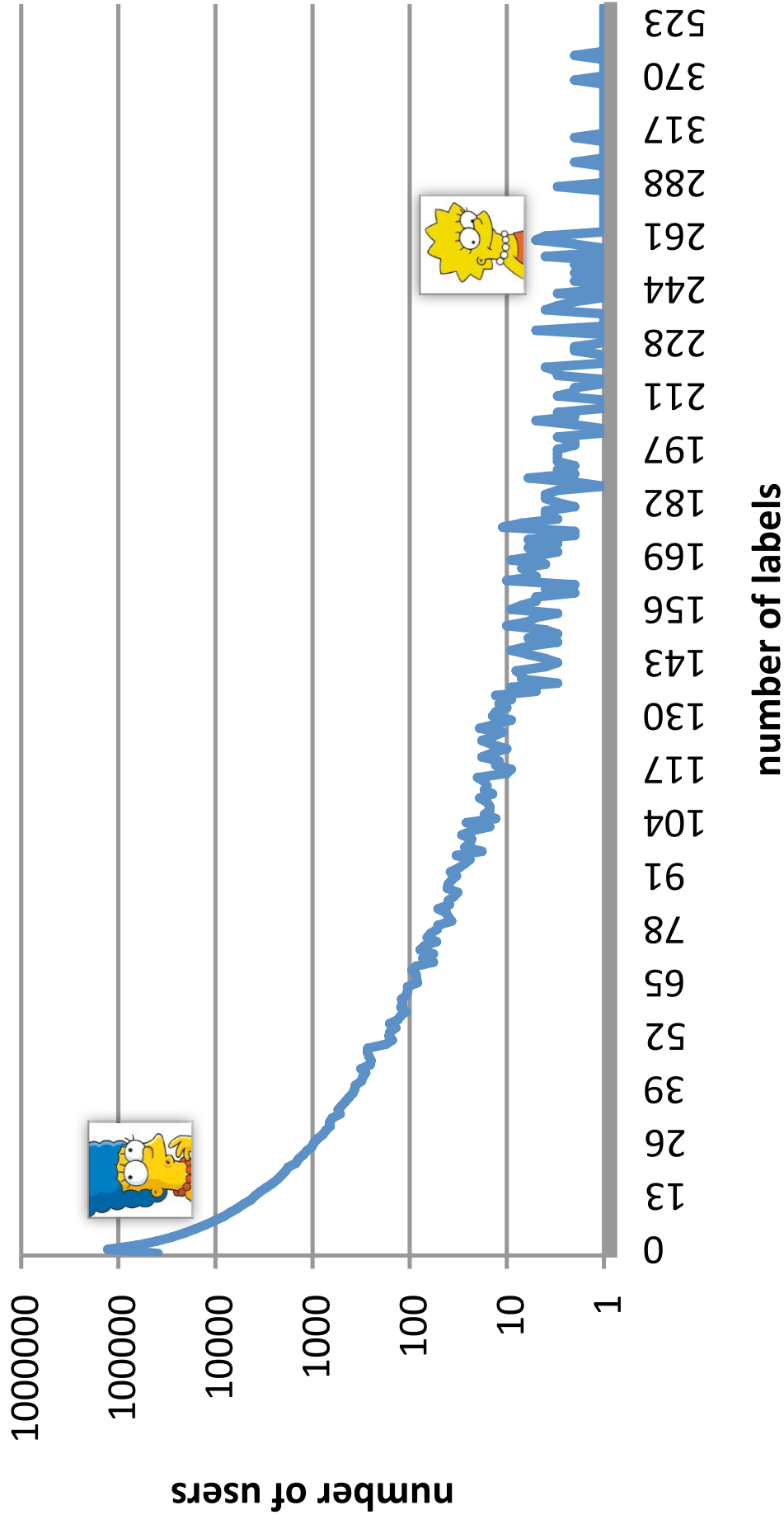# Spam classification results



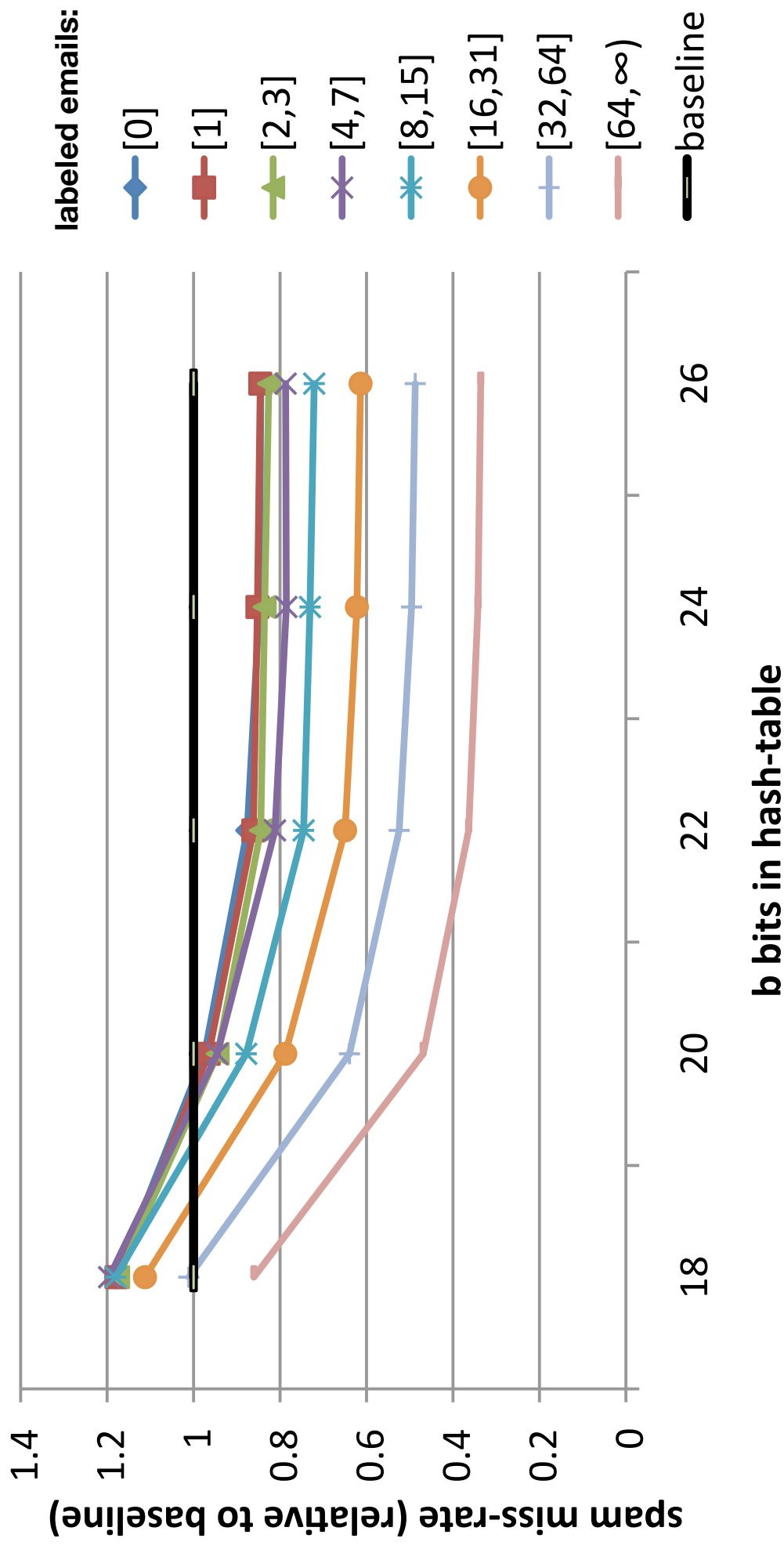spam miss-rate (relative to baseline)

b bits in hash-table

20 million emails, 400,000 users

- global-hashed
- personalized
- baseline

# Lazy users …

## Labeled emails per user



number of users

number of labels

# Results by user group



labeled emails:
- [0]
- [1]
- [2,3]
- [4,7]
- [8,15]
- [16,31]
- [32,64]
- [64,∞)
- baseline

spam miss-rate (relative to baseline)

b bits in hash-table

# Results by user group



labeled emails:
- [0]
- [1]
- [2,3]
- [4,7]
- [8,15]
- [16,31]
- [32,64]
- [64,∞)
- baseline

spam miss-rate (relative to baseline)

b bits in hash-table

Carnegie Mellon University

# Approximate String Matches

- General idea

$$k(x, x') = \sum_{w \in x} \sum_{w' \in x'} \kappa(w, w') \text{ for } |w - w'| \leq \delta$$

Carnegie
Carneg1e
0arnegie
Canegie
Carn3gie

**catch all**
**with wildcards**

Carnegie
Carneg*e
*arnegie
Ca*negie
Carn*gie

- Map into fragments: dog -> (*og, d*g, do*)
- Hash fragments and weigh them based on mismatch amount
- Exponential in number of mismatches. Not alphabet size.

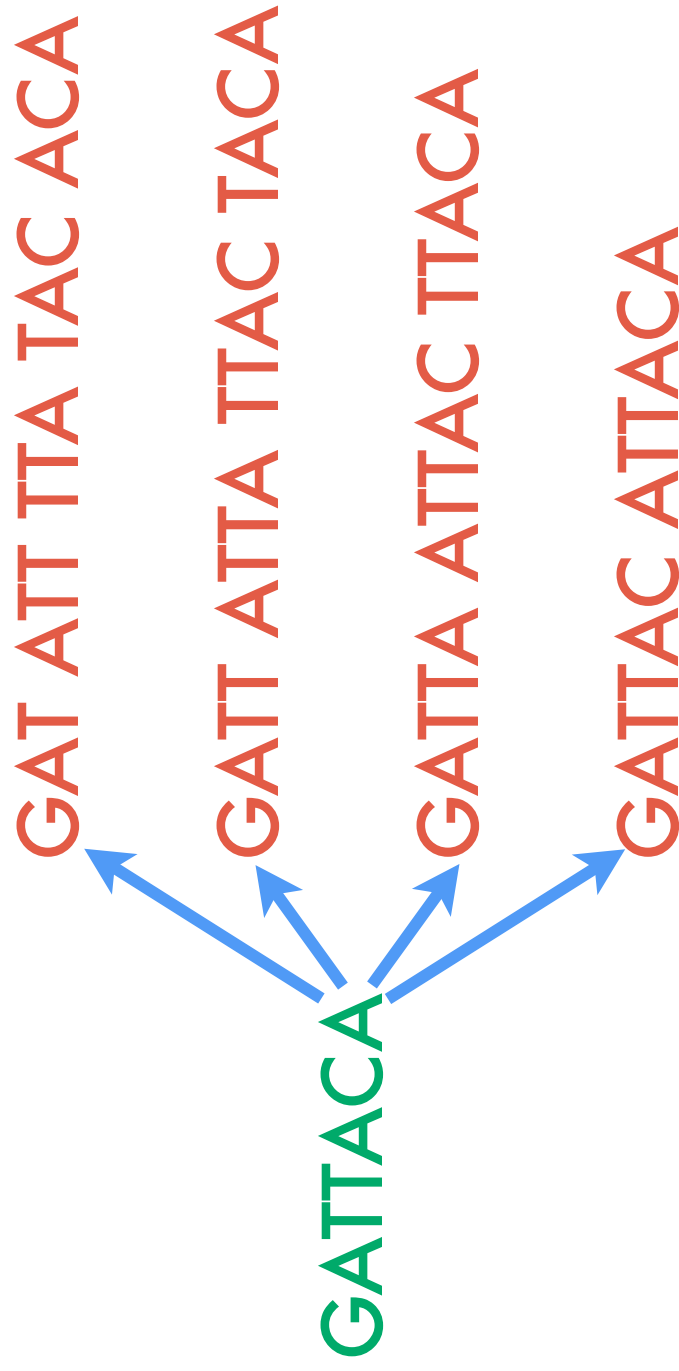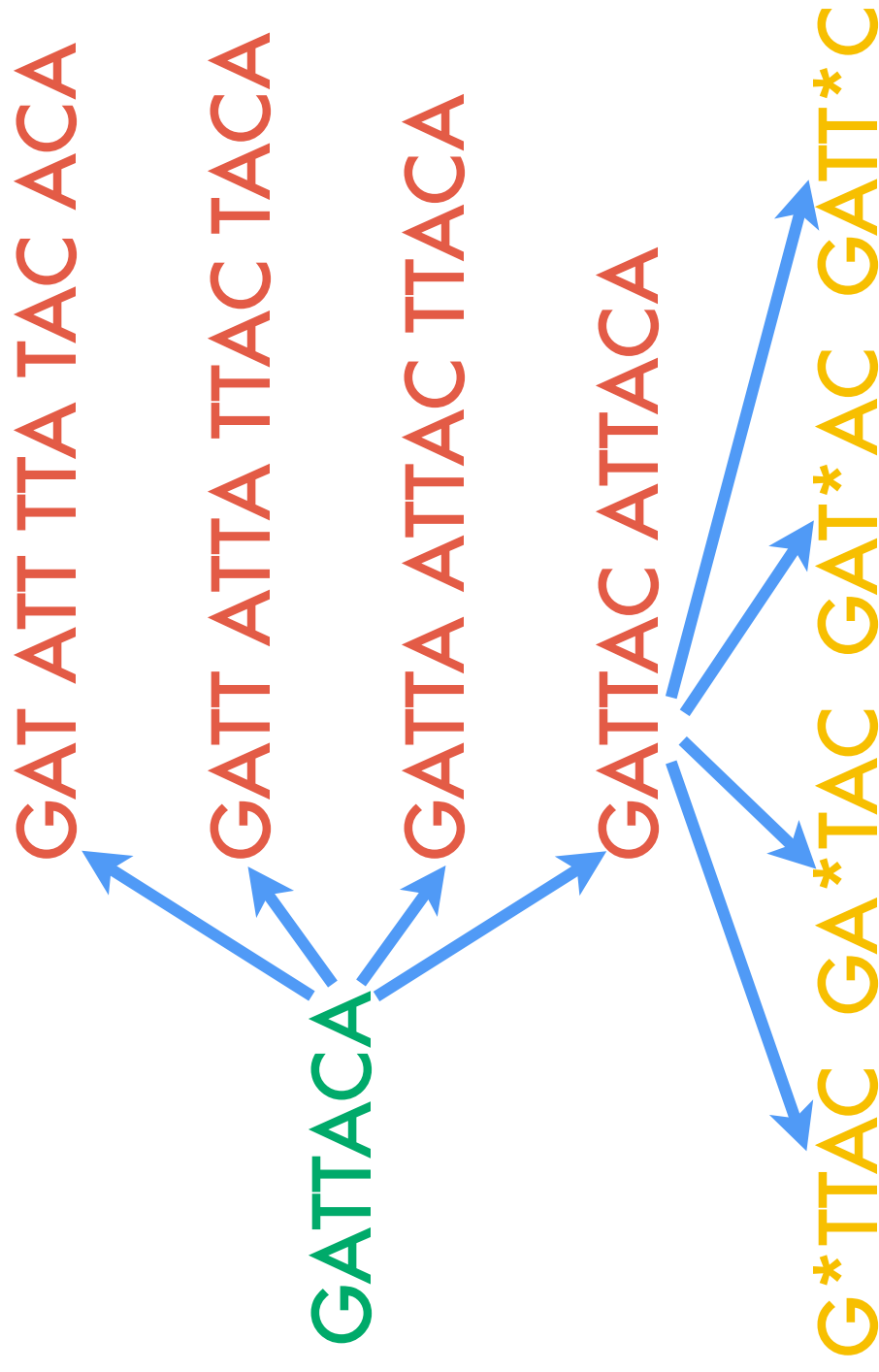# Fast String Kernels

- Example - DNA sequence
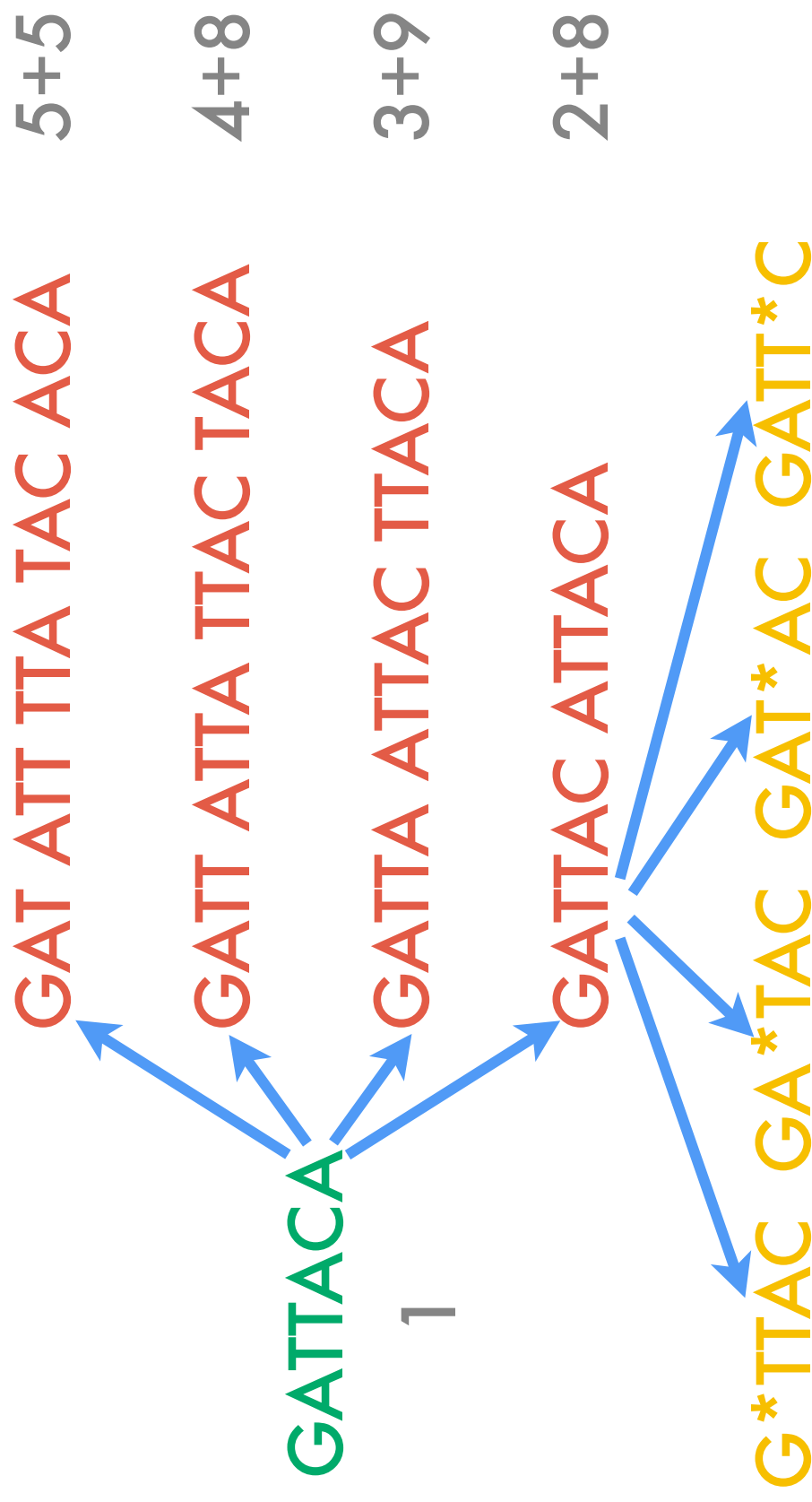
GATTACA

# Fast String Kernels

- Example - DNA sequence

GATTACA

GAT ATT TTA TAC ACA

GATT ATTA TTAC TACA

GATTA ATTAC TTACA

GATTAC ATTACA

# Fast String Kernels

- Example - DNA sequence

GATTACA

GAT ATT TTA TAC ACA

GATT ATTA TTAC TACA

GATTA ATTAC TTACA

GATTAC ATTACA

G*TTAC GA*TAC GAT*AC GATT*C

# Fast String Kernels

- Example - DNA sequence

GATTACA

1

GAT ATT TTA TAC ACA    5+5

GATT ATTA TTAC TACA    4+8

GATTA ATTAC TTACA    3+9

GATTAC ATTACA    2+8

G*TTAC GA*TAC GAT*AC GATT*C

# Fast String Kernels

- Example - DNA sequence

45

GATTACA

GAT ATT TTA TAC ACA    5+5

GATT ATTA TTAC TACA    4+8

GATTA ATTAC TTACA    3+9

GATTAC ATTACA    2+8

1

G*TTAC GA*TAC GAT*AC GATT*C

# Fast String Kernels

- Store coefficients explicitly (complicated)
- Use hash kernel to update counts (trivial)

GATTACA

GAT ATT TTA TAC ACA

GATT ATTA TTAC TACA

GATTA ATTAC TTACA

GATTAC ATTACA

G*TTAC  GA*TAC  GAT*AC  GATT*C