БПМИ-191: Абрамов Арсений

# MULTIPLYING MATRICES WITHOUT MULTIPLYING

Москва, 2022

- Матричное умножение лежит в основе многих алгоритмов и возникает почти всюду в ТИ и МО.

- Честное матричное умножение имеет кубическую сложность (плохо).

- Approximating Matrices Multiplication (AMM) ищет баланс между точностью и эффективностью по времени / памяти.

- Активно разрабатываются новые методы AMM, обладающие ещё большей эффективностью. Много методов AMM работают при допущениях.

Линейный оператор

$$A \in R^{N \times D}, B \in R^{D \times M}, N >> D \geq M$$

Data

$$V_A, V_B \in R^{D \times d}, d << D$$

sparse

$$AB \sim (AV_A)(V_B^T B)$$

(Если коротко, то проблема сводится к честному произведению в меньшей размерности через линейные преобразования A и B.)

# Assumptions:

Given a matrix $A \in \mathbb{R}^{m \times n}$, let us define:

ядро

- $A$ is a **fat matrix** if $m \leq n$ and $\mathbf{null}(A^T) = \{0\}$

- Matrices are tall

- $A$ is a **tall matrix** is $m \geq n$ and $\mathbf{range}(A) = \mathbb{R}^n$

образ

- Matrices are relatively dense

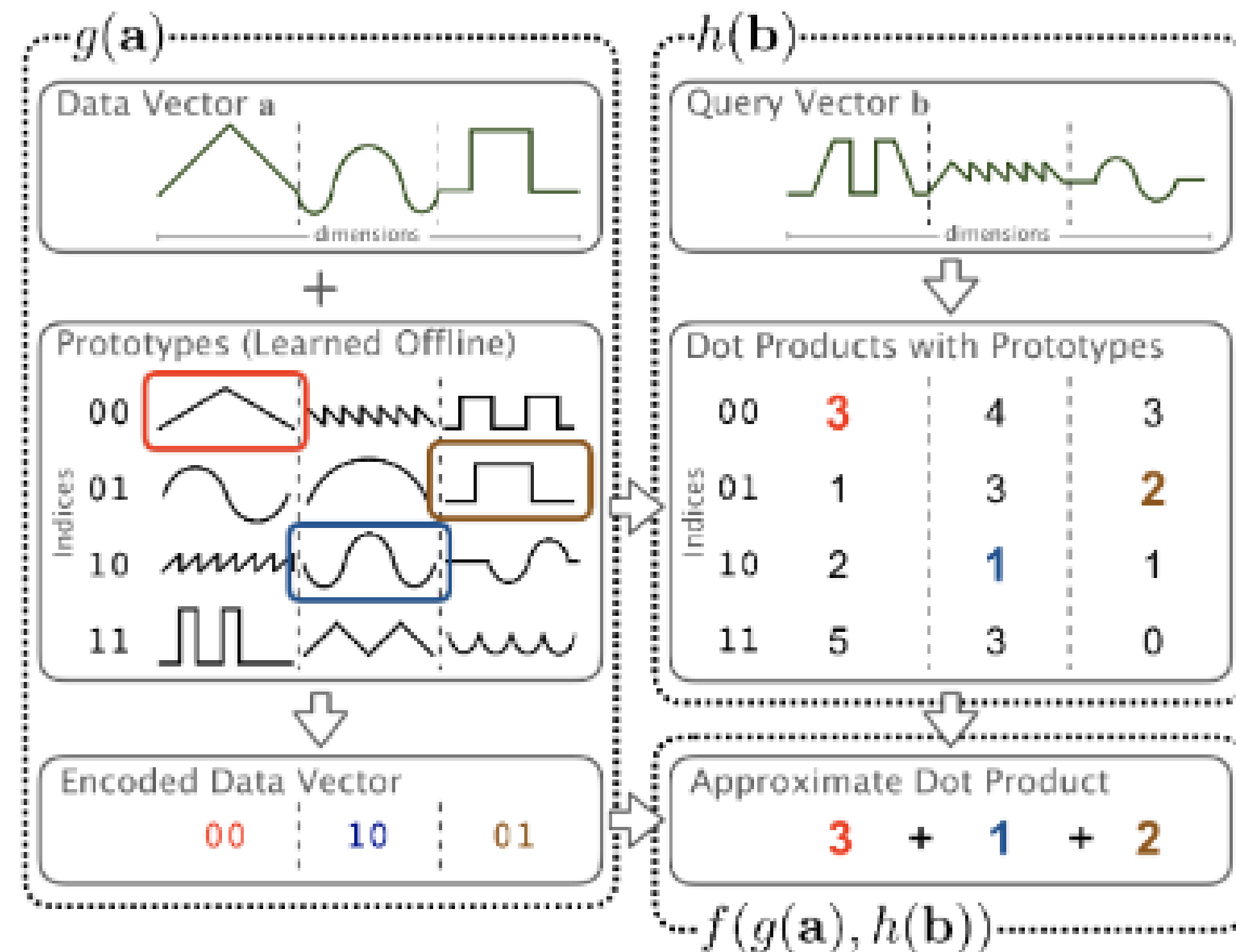- Resident in a single machine's memory

- **Exists a training set Ã.**

MADDNESS employs a nonlinear preprocessing function and reduces the problem to table lookups.

If B is *known ahead of time,* MADDNESS **requires no multiply-add operations.**

# Product Quantization



1. **Prototype Learning** - In an initial, offline training phase, cluster the rows of $A$ (or a training set $\bar{A}$) using K-means to create prototypes. A separate K-means is run in each of $C$ disjoint subspaces to produce $C$ sets of $K$ prototypes.

2. **Encoding Function**, $g(a)$ - Determine the most similar prototype to $a$ in each subspace. Store these assignments as integer indices using $C \log_2(K)$ bits.

3. **Table Construction**, $h(B)$ - Precompute the dot products between $b$ and each prototype in each subspace. Store these partial dot products in $C$ lookup tables of size $K$.

4. **Aggregation**, $f(\cdot, \cdot)$ - Use the indices and tables to *lookup* the estimated partial $a^\top b$ in each subspace, then sum the results across all $C$ subspaces.

# Hash function family

---
**Algorithm 1** MADDNESSHASH
---

1: **Input:** vector $x$, split indices $j^1, \ldots, j^4$, split thresholds $v^1, \ldots, v^4$
2: $i \leftarrow 1$      // node index within level of tree
3: **for** $t \leftarrow 1$ **to** $4$ **do**
4:    $v \leftarrow v_i^t$ // lookup split threshold for node $i$ at level $t$
5:    $b \leftarrow x_{j^t} \geq v$ ? $1 : 0$      // above split threshold?
6:    $i \leftarrow 2i - 1 + b$      // assign to left or right child
7: **end for**
8: **return** $i$

---

# Learning the hash function parameters

**Algorithm 2** Adding The Next Level to the Hashing Tree

1: **Input:** buckets $\mathcal{B}_1^{t-1}, \ldots, \mathcal{B}_{2^{t-1}}^{t-1}$, training matrix $\tilde{A}$

   // greedily choose next split index and thresholds

2: $\hat{\mathcal{J}} \leftarrow \texttt{heuristic\_select\_idxs}(\mathcal{B}_1^{t-1}, \ldots, \mathcal{B}_{2^{t-1}}^{t-1})$

3: $l^{min}, j^{min}, v^{min} \leftarrow \infty, \text{NaN}, \text{NaN}$

4: **for** $j \in \hat{\mathcal{J}}$ **do**

5:    $l \leftarrow 0$              // initialize loss for this index to 0

6:    $v \leftarrow [\,]$            // empty list of split thresholds

7:    **for** $i \leftarrow 1$ **to** $2^{t-1}$ **do**

8:       $v_i, l_i \leftarrow \texttt{optimal\_split\_threshold}(j, \mathcal{B}_i^{t-1})$

9:       $\texttt{append}(v, v_i)$   // append threshold for bucket $i$

10:      $l \leftarrow l + l_i$     // accumulate loss from bucket $i$

11:    **end for**

12:    **if** $l < l^{min}$ **then**

13:      $l^{min} \leftarrow l, j^{min} \leftarrow j, v^{min} \leftarrow v$  // new best split

14:    **end if**

15: **end for**

   // create new buckets using chosen split

16: $\mathcal{B} \leftarrow [\,]$

17: **for** $i \leftarrow 1$ **to** $2^{t-1}$ **do**

18:    $\mathcal{B}_{below}, \mathcal{B}_{above} \leftarrow \texttt{apply\_split}(v_i^{min}, \mathcal{B}_i^{t-1})$

19:    $\texttt{append}(\mathcal{B}, \mathcal{B}_{below})$

20:    $\texttt{append}(\mathcal{B}, \mathcal{B}_{above})$

21: **end for**

22: **return** $\mathcal{B}, l^{min}, j^{min}, v^{min}$

$$\mathcal{L}(j, \mathcal{B}) \triangleq \sum_{x \in \mathcal{B}} \left( x_j - \frac{1}{|\mathcal{B}|} \sum_{x' \in \mathcal{B}} x'_j \right)^2$$

$$\mathcal{L}(\mathcal{B}) \triangleq \sum_j \mathcal{L}(j, \mathcal{B}).$$

# Optimizing prototypes

$$P \triangleq (G^\top G + \lambda I)^{-1} G^\top \tilde{A}.$$

$P \in \mathbb{R}^{KC \times D}$

Param = 1
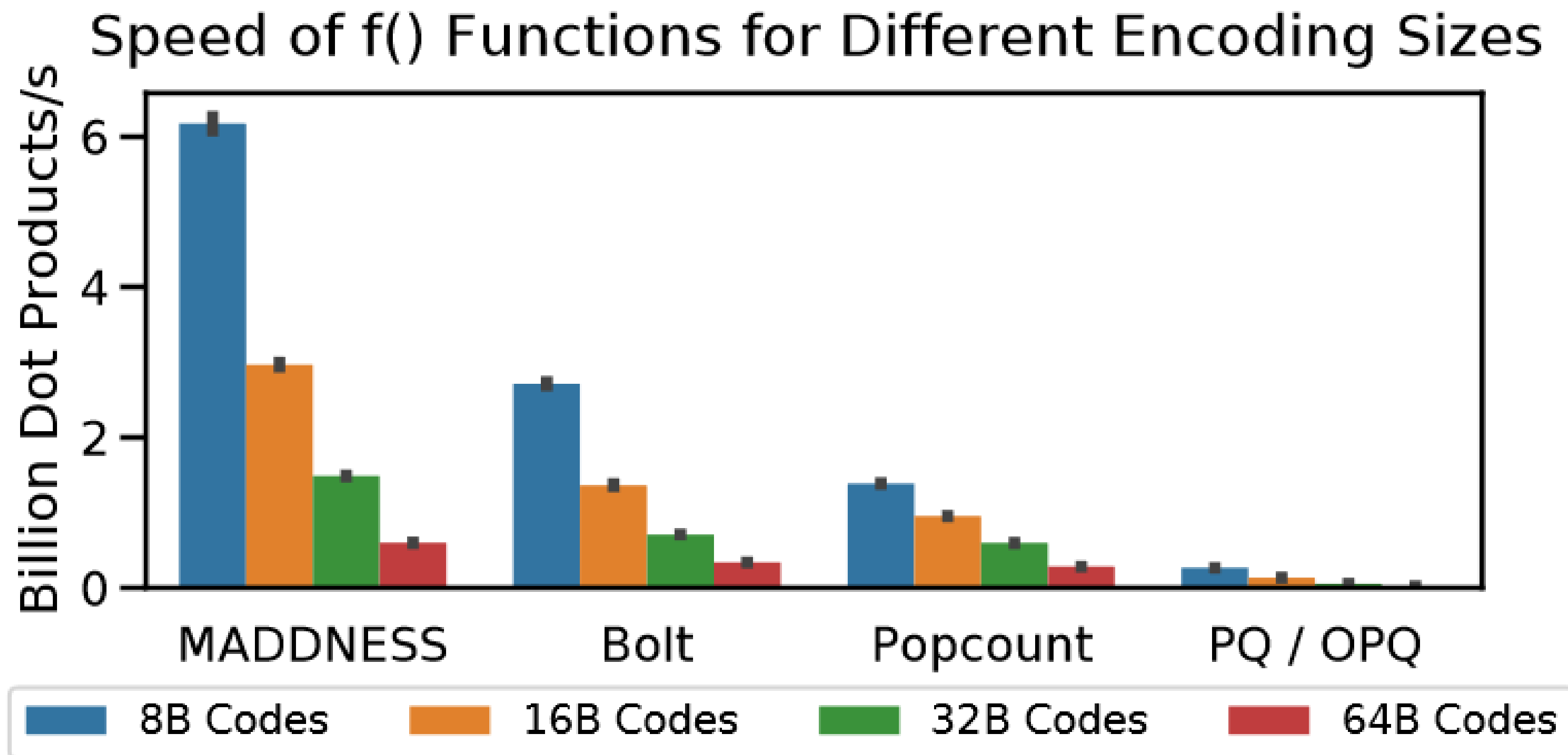
$\tilde{A} \approx \bar{G} P$

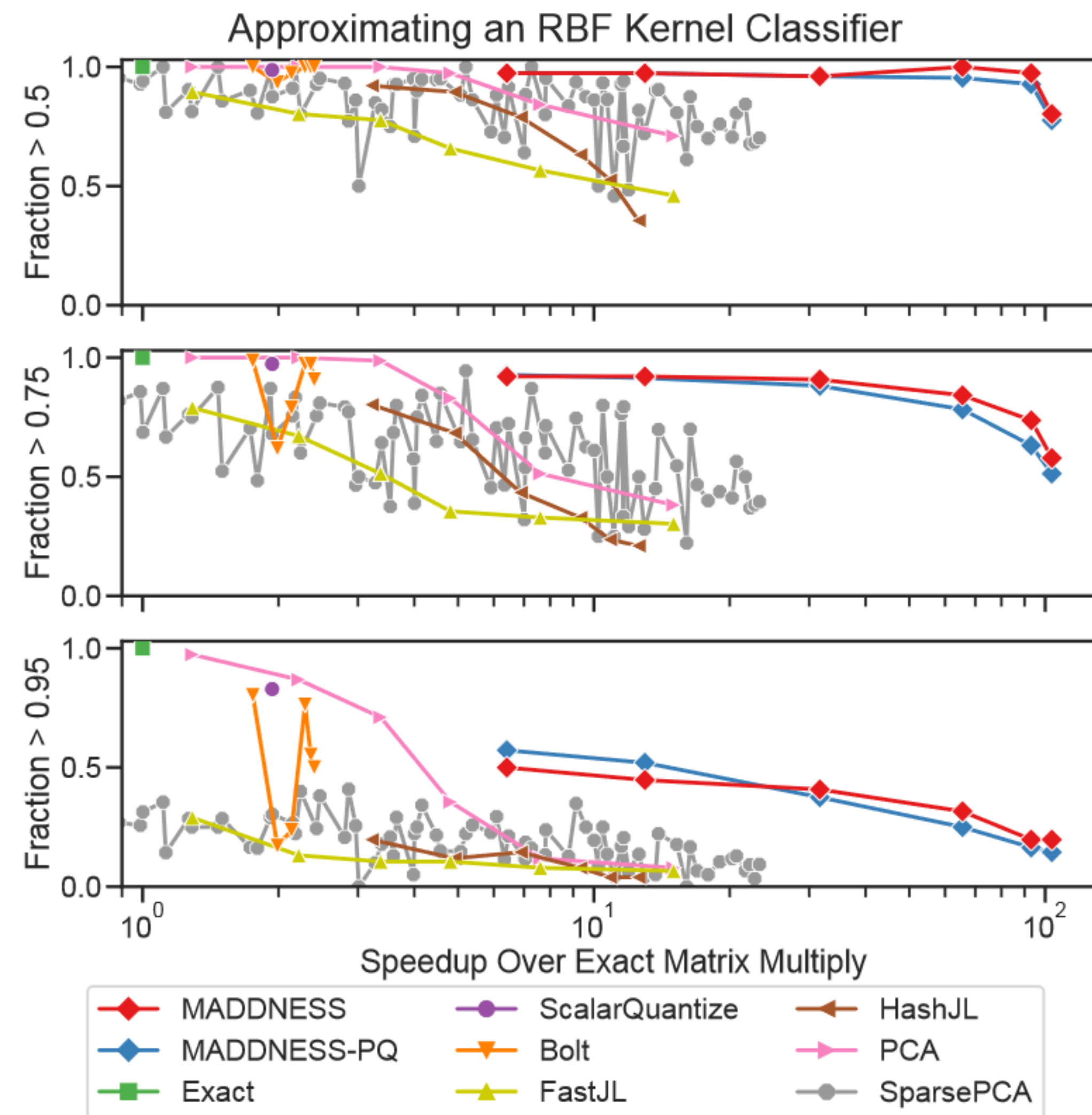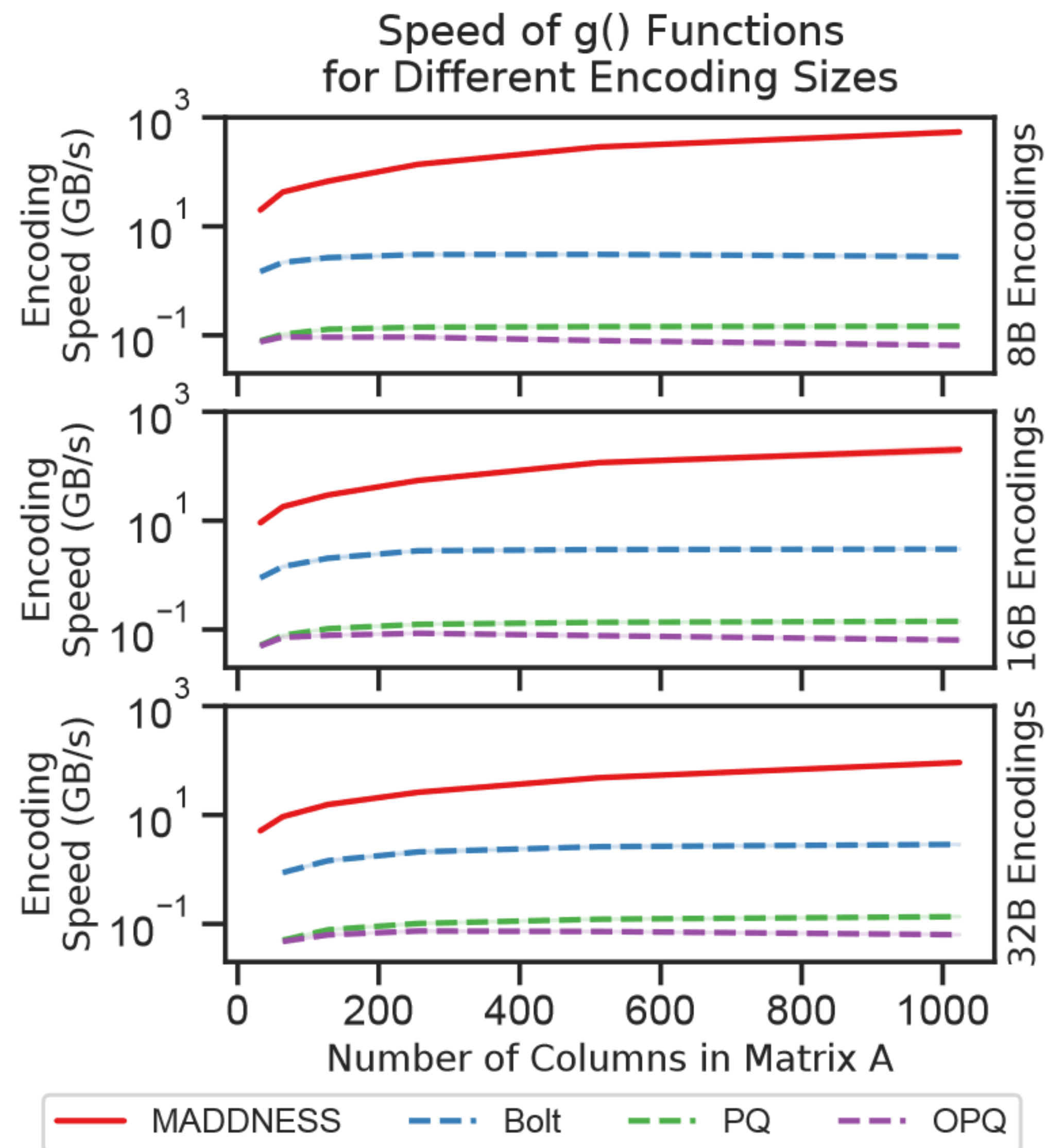## Fast 8-bit aggregation

Let $\boldsymbol{T} \in \mathbb{R}^{M \times C \times K}$ be the tensor of lookup tables for all $M$ columns of $\boldsymbol{B}$. Given the encodings $\boldsymbol{G}$, the function $f(\cdot, \cdot)$ is defined as

$$f(g(\boldsymbol{A}), h(\boldsymbol{B}))_{n,m} \triangleq \sum_{c=1}^{C} \boldsymbol{T}_{m,c,k}, \; k = g^{(c)}(\boldsymbol{a}_n). \quad (9)$$

~~Addition~~ —————————————→ Averaging

Speed of f() Functions for Different Encoding Sizes

$$AB \approx (AV_A)(V_B^\top B).$$

Learning Hash Function

# https://smarturl.it/Maddness

# ИСТОЧНИКИ

(Источник)

**Источник**

ИСТОЧНИК

- *Davis Blalock, John Guttag:* **Multiplying Matrices Without Multiplying**
  **URL:** https://arxiv.org/pdf/2106.10860.pdf

НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ