

Inner Product Similarity Join

Thomas D. Ahle, Rasmus Pagh, Ilya Razenshteyn, Francesco Silvestri



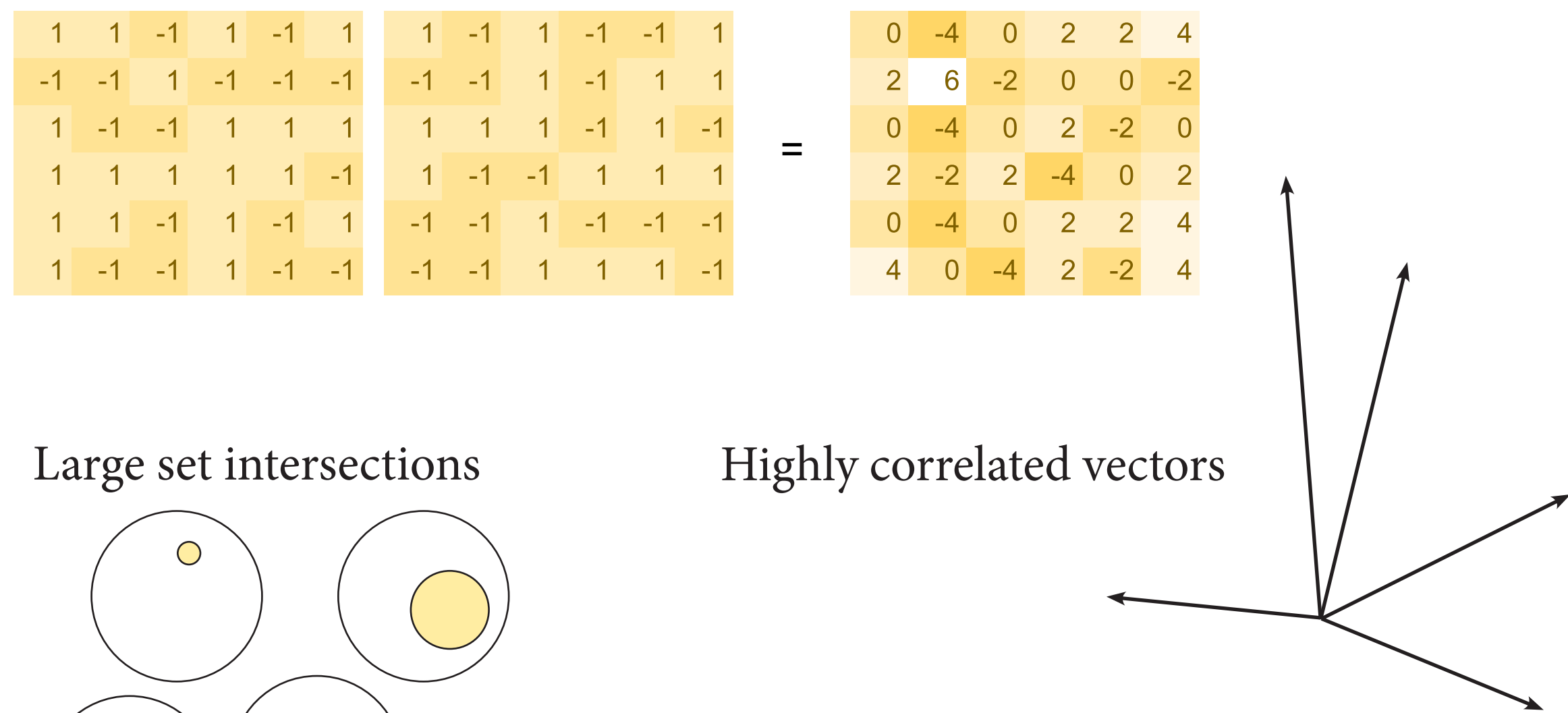
European Research Council
Established by the European Commission
Supporting top researchers
from anywhere in the world



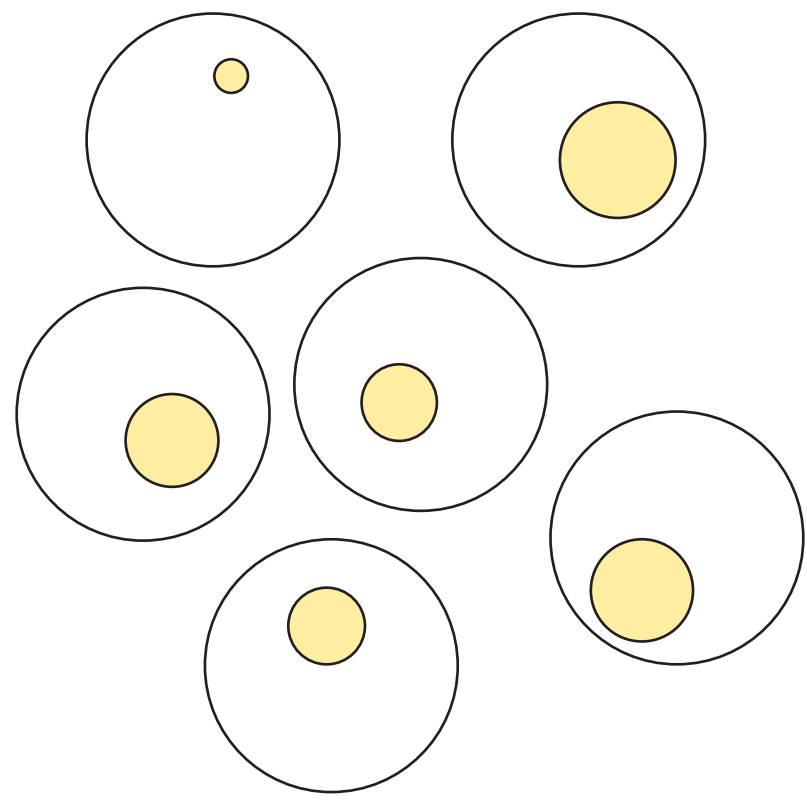
Problem definition

Exact IPS Join	c, s -Approx. Signed (c -join)	c -Approx. Unsigned
Given $P, Q \subseteq \mathbb{R}^n$ find all $x \in P, y \in Q$ st. $x^T y > s$	Given $P, Q \subseteq \mathbb{R}^n$ find $x \in P, y \in Q$ st. $x^T y > cs$	Given $P, Q \subseteq \mathbb{R}^n$ find $x \in P, y \in Q$ st. $ x^T y > cs$
	when it is guaranteed that a pair exists with inner product at least $s > 0$	when it is guaranteed that a pair exists with <u>absolute</u> inner product at least $s > cs$

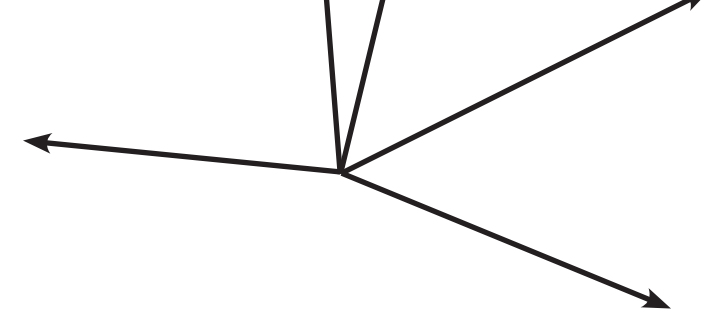
Large entries in matrix product



Large set intersections



Highly correlated vectors



Polynomial evaluations

$$4x^2 + 3xy - 2y^2 + 5 \quad x, y = 3, 4$$

$$8x^2 + 6xy + y^2 - 1 \quad x, y = -4, -5$$

$$x^2 - 20y^2 + 5 \quad x, y = 0, 0$$

Upper bounds - Exact IPS join

- Probabilistic Polynomials: $n^{2-1/O(c \log 2c)}$ for $P, Q \subseteq \{0, 1\}^{\log n}$ [Alman, Williams, FOCS 15]

Upper bounds - Approximate

- Data Dependent LSH: $n^{1+(1-s)/(1+s-2cs)}$ [Andoni, Razenshteyn STOC 15]
- Matrix multiplication: $n^{4/(3-\log s/\log cs)}$ [Karpka, Kaski, Kohonen, SODA 16]
- Sketching: For every $2 \leq \kappa \leq \infty$ there exists a distribution over $O(n^{1-2/\kappa}) \times n$ matrices Π such that for every $x \in \mathbb{R}^n$ one has: $\Pr_n[(1-c)\|x\|_\kappa \leq \|\Pi x\|_\infty \leq (1+c)\|x\|_\kappa] \geq 0.99$ [Andoni 10]

Main Theorem

There are no subquadratic algorithms for any of the following problems:

- Signed c -join for $c > 0$
- Unsigned c -join for $c = e^{-o(\sqrt{\log n / \log \log n})}$
- c -join of $P, Q \subseteq \{0, 1\}^d$ for $c = 1 - o(1)$
- Unsigned c -join for $\log(cs)/\log(s) = 1 - o(1/\sqrt{\log n})$
- c -join of $P, Q \subseteq \{0, 1\}^d$ for $\log(cs)/\log(s) = 1 - o(1/\log n)$



Created by Samu Parra from the Noun Project

The curse of dimensionality

Hardness in P "Understand hardness by reducing to well known hard problems"

- 3SUM: Given a set S of n integers, are there $a, b, c \in S$ with $a+b+c = 0$?
- Orthogonal vectors: Given a set S of n vectors in $\{0, 1\}^d$, for $d = O(\log n)$ are there $x, y \in S$ with $x^T y = 0$?
- All pairs shortest paths (APSP): given a weighted graph, find the distance between every two nodes.

Orthogonal Vectors Conjecture

No algorithm solves OV in time $O(n^{2-\epsilon})$ for $\epsilon > 0$

Proof Idea

1. We study embeddings $f, g: \{0, 1\}^{d_1} \rightarrow \mathbb{R}^{d_2}$ so for every $x, y \in \{0, 1\}^{d_1}$ $|f(x)^T g(y)| \leq cs$, if $x^T y \geq 1$ and $|f(x)^T g(y)| \geq s$, if $x^T y = 0$. "Make an embedding"
2. Make sure d_2 (and the embedding time) is $n^{o(1)}$. "Apply the embedding"
3. You've shown that unless OVP is false, there is no algorithm for unsigned (s, cs) join running in time $d^{o(1)} n^{2-\epsilon}$. "Solve OVP"



Simple embedding: showing $c > 0$ is impossible for the unsigned join

Take, coordinate wise

$$f(0) := (1, -1, -1) \quad g(0) := (1, 1, -1)$$

$$f(1) := (1, 1, 1) \quad g(1) := (-1, -1, -1)$$

Then

$$f(1)^T g(1) = -3$$

$$f(0)^T g(1) = f(1)^T g(0) = f(0)^T g(0) = 1$$

For vectors take

$$f(x) = f(x_1) \dots f(x_n) \quad 1^{d-4}$$

$$g(x) = g(x_1) \dots g(x_n) \quad (-1)^{d-4}$$

Such that

$$f(x)^T g(y) \leq 0, \quad \text{if } x^T y \geq 1$$

$$f(x)^T g(y) = 4, \quad \text{if } x^T y = 0$$

For the other results, we use polynomial embeddings ~ [Valiant, FOCS 12]

$\{-1, 1\}$ Hardness by Chebyshev Polynomials

$$T_n(x) = \cos(n \arccos(x))$$

$$|T_n(x)| \leq 1 \quad \text{for } -1 \leq x \leq 1$$

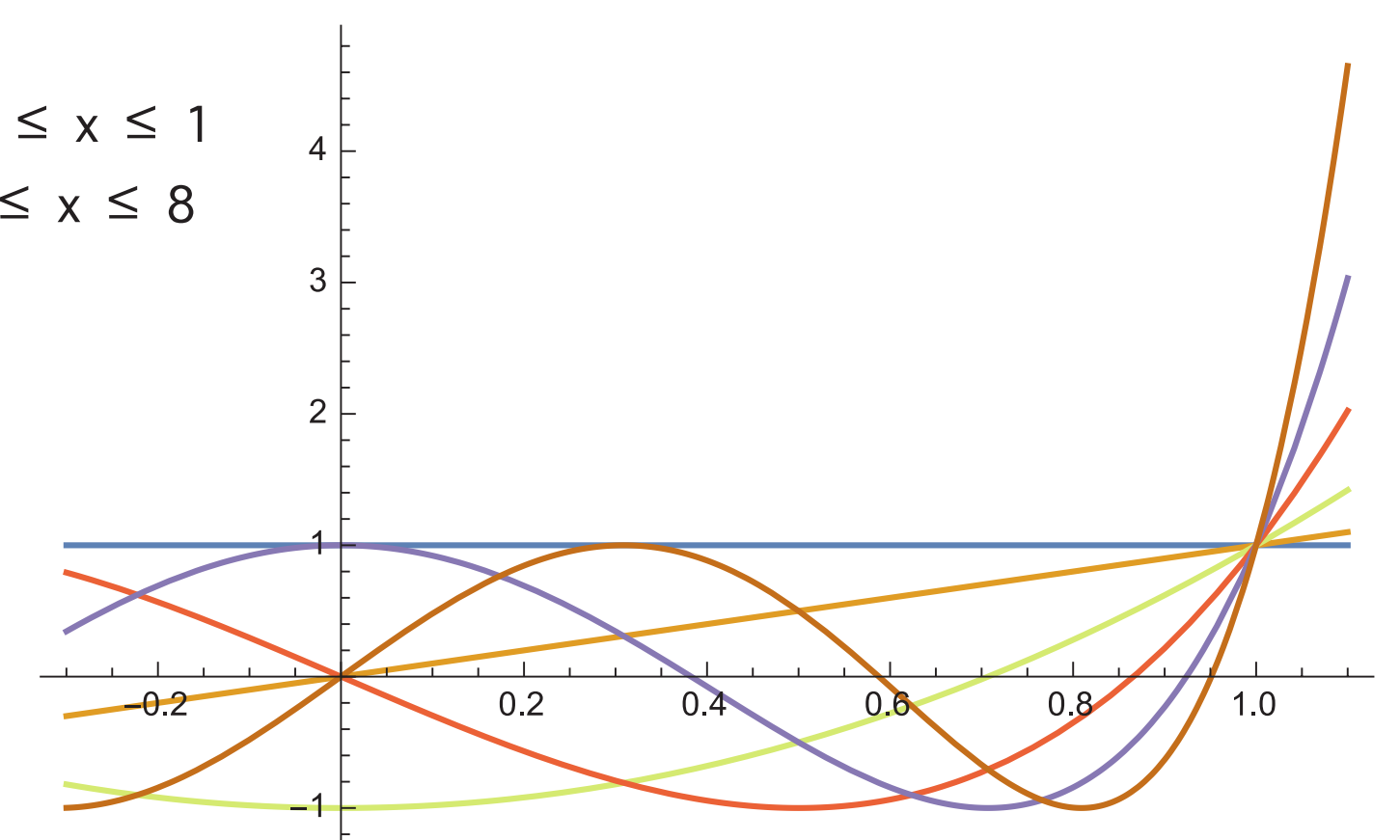
$$T_n(x) \geq e^{n \sqrt{x-1}} \quad \text{for } 1 \leq x \leq 8$$

$$T_0(x) = 1$$

$$T_1(x) = x$$

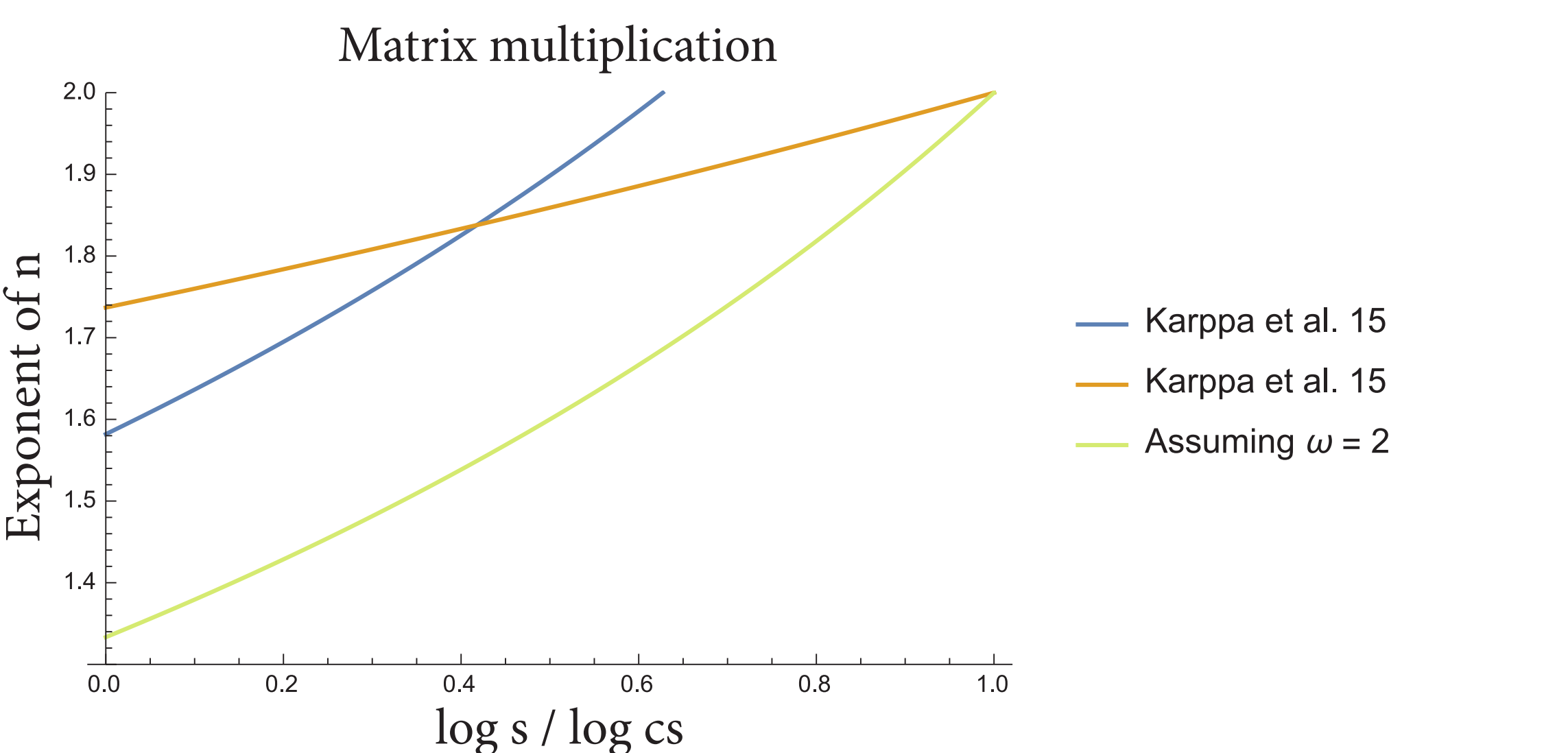
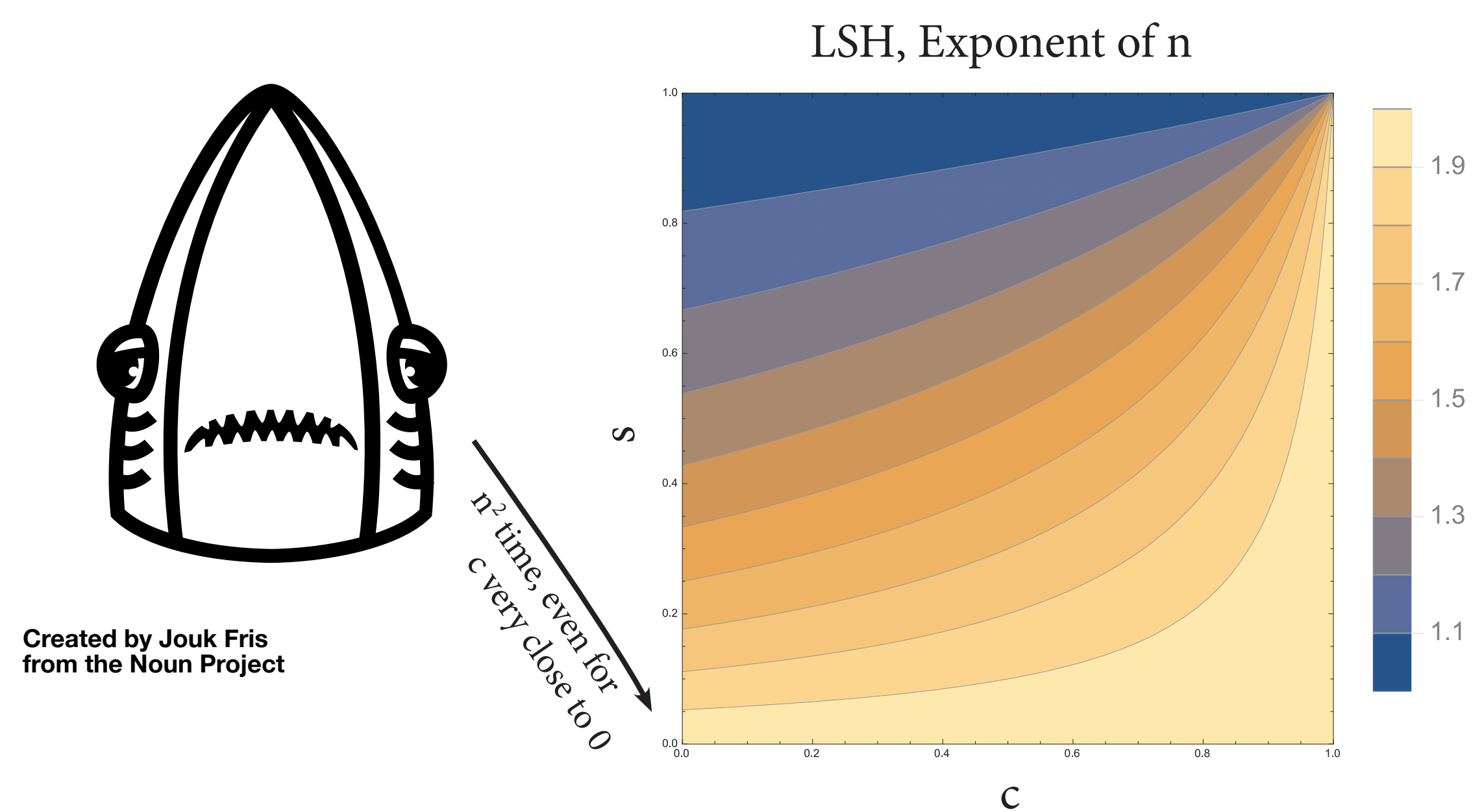
$$T_2(x) = 2x^2 - 1$$

$$T_3(x) = 4x^3 - 3x$$



$\{0, 1\}$ Hardness by truncated product construction

$$T(x) = (1-x_1 y_1) \dots (1-x_k y_k) + (1-x_{k+1} y_{k+1}) \dots (1-x_{2k} y_{2k}) + \dots$$



Other results	Future work
<ul style="list-style-type: none"> Hardness for datastructures, via reduction through asymmetric OVP. Lower bound for LSH, "p1 - p2 = O(sqrt(s))". Elimination of asymmetry via error correcting codes. 	<ul style="list-style-type: none"> Get tighter bounds. Can we disprove subquadratic running time for logs/logcs = 1-eps? Prove or disprove separation between $\{-1, 1\}$ and $\{0, 1\}$ cases. Reconcile LSH and Matrix Multiplication methods.