

# Fast (and yet correct!) Tensor Processing

**Presented by: Ana Sá Oliveira, Edgar Araújo, Gabriel Paiva**

Supervised by: José Nuno Oliveira

University of Minho, Informatics Department

March 2026



# Overview

Context and Motivation: Tensor Processing

Problem: Untyped Linear Algebra

Proposed Solution: Category Theory

Case Study: AXPY

Literature Review

Future Work & Schedule

# Context and Motivation: Tensor Processing

# Tensor Processing

Tensor processing is an advanced form of **Basic Linear Algebra (BLAS)**.

- ▶ Foundation of numeric programming
- ▶ Used in computational fluid dynamics and large-scale simulations (tsunamis, epidemics)
- ▶ Foundation of **Machine Learning**, in particular, Large Language (Transformer) Models (**LLMs**)
- ▶ Important for Probabilistic Programming



ChatGPT

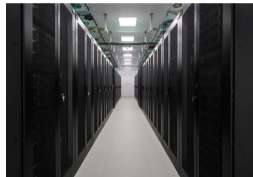


Gemini

# **Problem: Untyped Linear Algebra**

# Prioritizing Performance over Safety

- ▶ Large-scale tensor processing relies on **HPC** and massively parallel architectures, prioritizing **performance** over safety.
- ▶ **Manual coding** predominates to ensure optimal performance.
- ▶ **Automated tools aren't trusted** to deliver both performance and correctness simultaneously.



# Untyped Linear Algebra

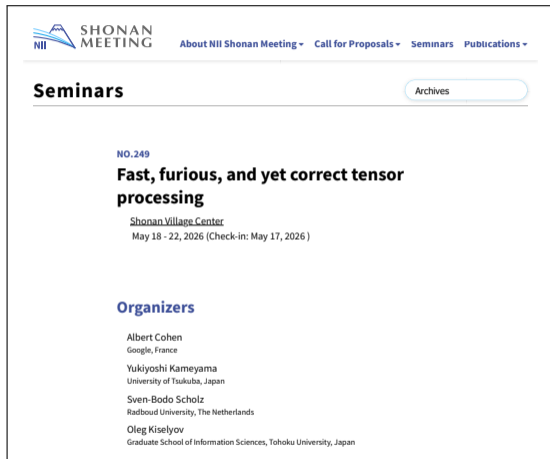
- ▶ Frameworks (e.g., PyTorch, TensorFlow) **lack static dimensional type checking**.
- ▶ *Shape mismatch* errors happen at runtime → **Wasted time and money**.



TensorFlow

# Current Approaches

- ▶ Rank-polymorphism
- ▶ Tensor comprehensions & transformations
- ▶ **Biproducts**
- ▶ APL array programming patterns
- ▶ Exterior algebra
- ▶ Shape calculi / size types



The screenshot shows the SHONAN MEETING website. At the top left is the NII logo and the text 'SHONAN MEETING'. To the right are navigation links: 'About NII Shonan Meeting - Call for Proposals - Seminars Publications -'. Below the navigation is a 'Seminars' section with an 'Archives' button. The featured seminar is NO.249, titled 'Fast, furious, and yet correct tensor processing', held at the Shonan Village Center from May 18 - 22, 2026. The organizers listed are Albert Cohen (Google, France), Yuki Yoshi Kameyama (University of Tsukuba, Japan), Sven-Bodo Scholz (Radboud University, The Netherlands), and Oleg Kiselyov (Graduate School of Information Sciences, Tohoku University, Japan).

*Meeting Goal:* Bring isolated communities together to stimulate convergence.

# Proposed Solution: Category Theory

# The Category of Matrices ( $\mathbf{Mat}_K$ )

In the category  $\mathbf{Mat}_K$ :

- ▶ **Objects:** Natural numbers ( $n, m \in \mathbb{N}$ ) representing dimensions.
- ▶ **Morphisms:** A morphism  $A$  from  $n$  to  $m$  is an  $m \times n$  matrix.

## The Arrow (Categorical)

$$\begin{array}{ccc} m & \xleftarrow{A} & n \\ \text{(Rows)} & & \text{(Columns)} \end{array}$$

## The Matrix (Concrete)

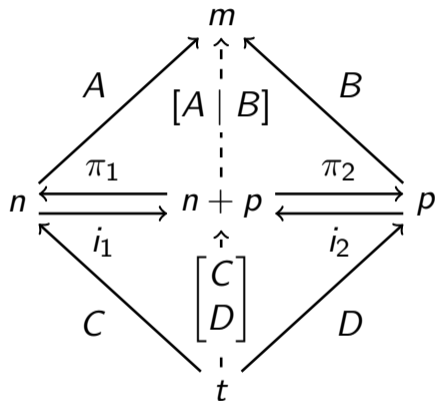
$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix}$$

## Biproducts: Formalizing Vectorization & Parallelism

We rely on the categorical notion of a **Biproduct** to provide a strict foundation for blocked linear algebra.

This structure allows us to formally specify:

- ▶ **Vectorization:** Packing data horizontally for optimized operations.
- ▶ **Parallelism:** Splitting computations in blocks to distribute workloads.

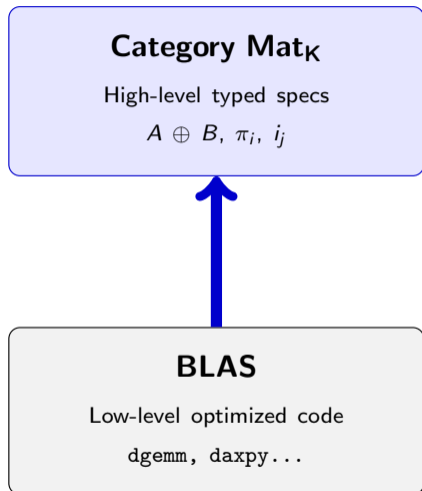


# Methodology: Reverse Engineering BLAS

## Reverse Engineering

via Categorical &  
Biproduct Laws

Statically verified type  
checking



# Case Study: AXPY

# The AXPY Algorithm

The name **AXPY** stands for the fundamental Linear Algebra operation:  $y = a \cdot x + y$ .

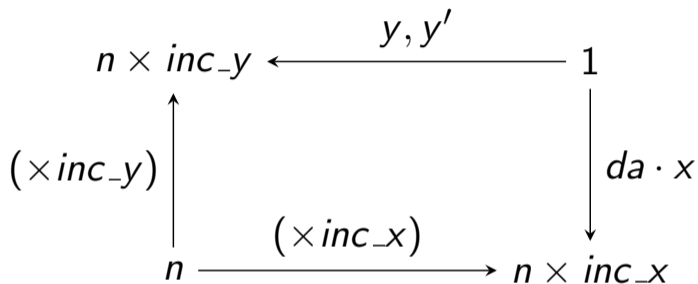
## Pseudo-code

$$y'[i \times inc\_y] = y[i \times inc\_y] + da \cdot x[i \times inc\_x]$$

## Visual Example (Vectors):

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} + 3 \cdot \begin{bmatrix} 4 \\ 5 \\ 6 \\ 3 \end{bmatrix} \quad (\text{considering strided access to } x)$$

## Type Diagram (Pointwise)



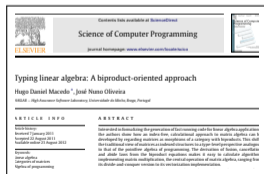
- ▶ **Vertical/Horizontal arrows:** Represent index scaling by strides.
- ▶ **Top/Right arrows:** Represent the data access and operations ( $y$  and  $x$ ).

# Literature Review

# Literature Review: Foundational Papers

Our review bridges formal theory and industry practice, drawing from foundational papers on biproduct-oriented linear algebra and state-of-the-art **HPC** technical literature:

## Typing linear algebra: A biproduct-oriented approach



The image shows the front page of the journal 'Science of Computer Programming'. At the top, it says 'Contents lists available at ScienceDirect' and 'Journal homepage: www.elsevier.com/locate/scp'. The title of the article is 'Typing linear algebra: A biproduct-oriented approach' by Hugo Daniel Macedo and José Nuno Oliveira. Below the title, it lists the authors' affiliations: 'HPC@FEUP - High-Performance Software Laboratory, Universidade do Porto, Porto, Portugal'. There are two columns: 'ARTICLE INFO' and 'ABSTRACT'. The abstract text is partially visible and discusses formalizing the generation of low-level code for linear algebra applications.

Contents lists available at ScienceDirect

Science of Computer Programming

Journal homepage: [www.elsevier.com/locate/scp](http://www.elsevier.com/locate/scp)

Typing linear algebra: A biproduct-oriented approach

Hugo Daniel Macedo<sup>a</sup>, José Nuno Oliveira

<sup>a</sup>HPC@FEUP - High-Performance Software Laboratory, Universidade do Porto, Porto, Portugal

---

**ARTICLE INFO**

**ABSTRACT**

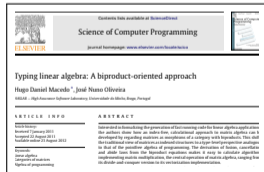
Formalizing the generation of low-level code for linear algebra applications has not been an end-to-end, industrial approach to matrix algebra, but is developed by regarding matrices as monomorphisms of a category with biproducts. This shift the traditional view of matrices as mathematical objects to larger, more expressive abstractions in that of the *prolinear algebra* of programming. The derivation of basic, canonical and stable linear-time, linear-branching algorithms to solve in solution algorithms implementing matrix multiplication, the creation of operations of matrix algebra, ranging from its divide-and-conquer, reduce to its vector-linear implementations.

# Literature Review: Foundational Papers

Our review bridges formal theory and industry practice, drawing from foundational papers on biproduct-oriented linear algebra and state-of-the-art **HPC** technical literature:

Typing linear algebra: A biproduct-oriented approach

The data cube as a typed linear algebra operator



Contents lists available at ScienceDirect  
Science of Computer Programming  
Journal homepage: [www.elsevier.com/locate/scp](http://www.elsevier.com/locate/scp)

Typing linear algebra: A biproduct-oriented approach

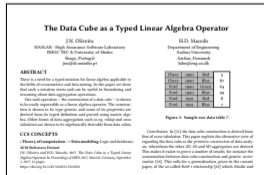
Hugo Daniel Macedo<sup>a</sup>, José Nuno Oliveira

<sup>a</sup>RIIS - High-Performance Software Laboratory, Universidade de Lisboa, Campo Real

**ARTICLE INFO**

**ABSTRACT**

Interested in formalizing the generation of fast-running code for linear algebra applications that involves dense, sparse or mixed forms, industrial approaches to matrix algebra can be developed by regarding matrices as monomials of a category with biproducts. This shift in the traditional view of matrices is combined with the algebraic linear programming strategies in that of the *procedural algebra of programming*. The derivation of linear, cancellative and additive Term Term (the Operational Equational Theory) is used to calculate algorithms implementing matrix multiplication, the evaluation of matrix algebra, ranging from its divide-and-conquer version to the vector-tensor implementation.



The Data Cube as a Typed Linear Algebra Operator

J.N. Oliveira  
RIIS/LAL - High-Performance Software Laboratory  
INESC TEC / University of Lisbon  
Ruiyi Peng  
juno@isec.pt

H.D. Macedo  
Department of Engineering  
Faculty University  
Lisboa, Portugal  
hdma@eng.ucp.li

**ABSTRACT**

There is a need for a typed calculus for linear algebra applicable to the fields of mathematics and data mining. In this paper we show that such a calculus exists and can be used to formalizing and reasoning about data-algebraic operations.

Our main operation – the construction of the data cube – is shown to be readily representable as a linear algebra operator. The construction is shown to be type-agnostic and even of the program, and defined from its typed definition and generalizing matrix algebra. Other forms of data aggregation such as *avg*, *min* and *max* calculation are shown to be algorithmically derivable from these others.

**CCS CONCEPTS**

**Theory of computation** → Data modeling; Logic and semantics; ACM Reference format

J.N. Oliveira and H.D. Macedo. 2017. The Data Cube as a Typed Linear Algebra Operator. In *Proceedings of the 2017 ACM SIGPLAN International Conference on Functional Programming (FPC '17)*, pages 1–11.

*https://doi.org/10.1145/3124210*

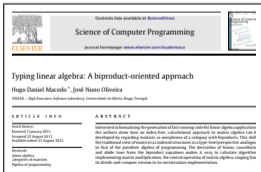
**Figure 1. Sample view data table 7.**

Cherry	Apple	Blue	0
Cherry	Apple	Blue	10
Apple	Apple	Green	14
Apple	Apple	Blue	14
Apple	Apple	Blue	14
Apple	Apple	Blue	14

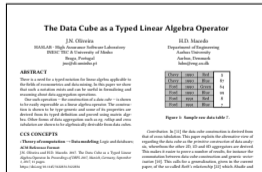
# Literature Review: Foundational Papers

Our review bridges formal theory and industry practice, drawing from foundational papers on biproduct-oriented linear algebra and state-of-the-art **HPC** technical literature:

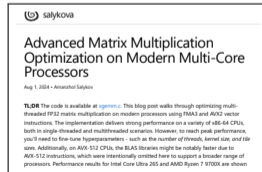
## Typing linear algebra: A biproduct-oriented approach



## The data cube as a typed linear algebra operator



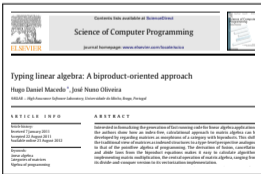
## Advanced Matrix Multiplication Optimization



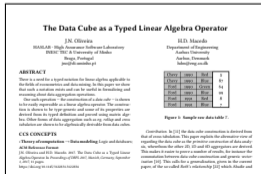
# Literature Review: Foundational Papers

Our review bridges formal theory and industry practice, drawing from foundational papers on biproduct-oriented linear algebra and state-of-the-art **HPC** technical literature:

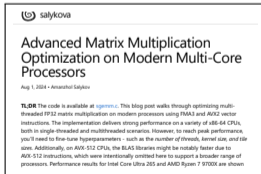
## Typing linear algebra: A biproduct-oriented approach



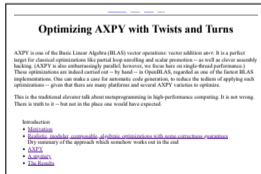
## The data cube as a typed linear algebra operator



## Advanced Matrix Multiplication Optimization

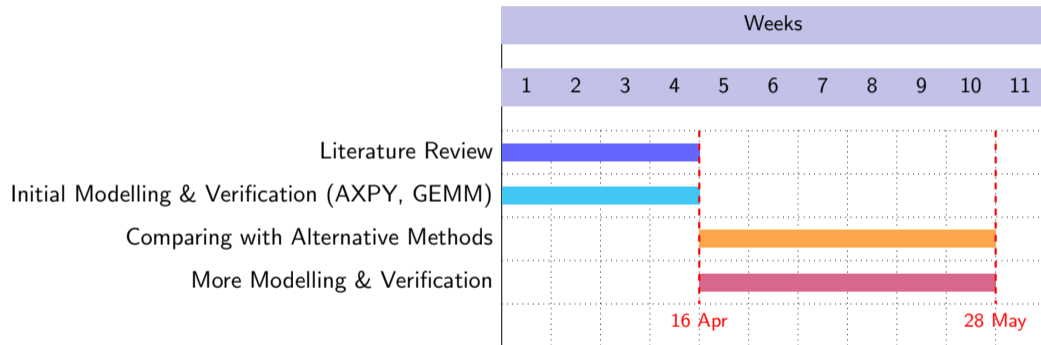


## Optimizing AXPY with Twists and Turns



# **Future Work & Schedule**

# Future Work & Schedule



# Fast (and yet correct!) Tensor Processing

**Presented by: Ana Sá Oliveira, Edgar Araújo, Gabriel Paiva**

Supervised by: José Nuno Oliveira

University of Minho, Informatics Department

March 2026

