EuroLLVM 2025 - Berlin

# MLIR Tensor Compiler

### design group & charter update

Rolf Morel & Renato Golin



# Upstream MLIR Tensor Compiler

- MLIR-based toolkit for ML compilers
  - Great value in developing this together upstream
- Main dialects:
  - linalg
  - tensor
  - memref
  - vector
  - bufferization
  - tosa



2

### MLIR organization – dialect groupings

- Originates with MLIR Governance efforts from last year
  - [RFC] MLIR Project Charter and Restructuring by Renato Golin, Stella Laurenzo, Chris Lattner, Alex Zinenko, Jacques Pienaar, Mehdi Amini
  - [Survey] MLIR Project Charter and Restructuring Survey by Renato
  - MLIR Organization & Charter by Renato, Stella, Alex, Jacques, Chris, and Andrzej Warzynski, Nicolas Vasilache, Mahesh Ravishankar



4

### MLIR organization – dialect groupings

- Originates with MLIR Governance efforts from last year
  - [RFC] MLIR Project Charter and Restructuring by Renato Golin, Stella Laurenzo, Chris Lattner, Alex Zinenko, Jacques Pienaar, Mehdi Amini
  - [Survey] MLIR Project Charter and Restructuring Survey by Renato
  - MLIR Organization & Charter by Renato, Stella, Alex, Jacques, Chris, and Andrzej Warzynski, Nicolas Vasilache, Mahesh Ravishankar
- Tensor Compiler dialect grouping
  - linalg, tensor, memref, vector, bufferization, tosa
  - Own community of stakeholders
  - In need of an overarching charter
  - Clear governance: upstream consensus

Linalg Connectivity

#### MLIR Tensor Compiler Design Group by Renato

Goal: charter across relevant dialects, interfaces, transforms, surrounding infrastructure



5

### MLIR Tensor Compiler Design Group by Renato

Goal: charter across relevant dialects, interfaces, transforms, surrounding infrastructure

### Scope

- linear algebra tensor semantics bufferization memref semantics vector semantics
- canonicalization (op aliasing within linalg) ingress & egress dialect design flows



### MLIR Tensor Compiler Design Group by Renato

Goal: charter across relevant dialects, interfaces, transforms, surrounding infrastructure

### Scope

- linear algebra tensor semantics bufferization memref semantics vector semantics
- canonicalization (op aliasing within linalg) ingress & egress dialect design flows

### **Documentation updates**

- Consolidate conflicting roadmaps
- Update outdated rationale docs
- Start an overarching charter



### MLIR Tensor Compiler Design Group by Renato

Goal: charter across relevant dialects, interfaces, transforms, surrounding infrastructure

### Scope

- linear algebra tensor semantics bufferization memref semantics vector semantics
- canonicalization (op aliasing within linalg) ingress & egress dialect design flows

### **Documentation updates**

- Consolidate conflicting roadmaps
- Update outdated rationale docs
- Start an overarching charter

### Infrastructure for distributed usage

- Devise a path for more flexibility for users (also across upstream projects)
- E.g., easier dialect extensions, canonicalized transform requirements, better coverage

8

### Tensor Compiler Design Group – Members

**O**<u>Volunteers</u> from active contributors and representative stakeholders

Alex Zinenko	Renato Golin	Jacques Pienaar	Matthias Springer
Brium	Intel	Google	Nvidia
Quinn Dawkins	Javed Absar	Jakub Kuderski	Suraj Sudhir
AMD	Qualcomm	AMD	Arm
Rolf Morel	Andrzej Warzynski	Diego Caballero	Kunwar Grover
Intel	Arm	Nvidia	AMD



### Tensor Compiler Design Group – Members

**O**<u>Volunteers</u> from active contributors and representative stakeholders

<i>Alex Zinenko</i>	<i>Renato Golin</i>	<i>Jacques Pienaar</i>	<i>Matthias Springer</i>
Brium	Intel	Google	Nvidia
Quinn Dawkins	Javed Absar	Jakub Kuderski	Suraj Sudhir
AMD	Qualcomm	AMD	Arm
Rolf Morel	Andrzej Warzynski	Diego Caballero	Kunwar Grover
Intel	Arm	Nvidia	AMD

\*In bold: MLIR Area Team, a distinct governance effort



Tensor Compiler Design Group

# Tensor Compiler Design Group – Members

<mark>)</mark> <u>Volunteers</u> fr	<b>Nodular</b> Democratizing AI Compute, Part 8: What about the MLIR compiler infrastructure?	
	A New Hope: Improved MLIR Governance	
<i>Alex Zinenko</i> Brium	The tensions have simmered for years—and they're deeply felt across the broader Atthias Springer LLVM and MLIR communities.	
Quinn Dawkins AMD	Fortunately, <b>there's a new hope</b> : LLVM is a meritocratic community with a long track record of aligning engineers—even when their companies are at war in the market. The MLIR community is filled with amazing engineers who have poured years of their hearts and souls into improving the project to work through these challenges, and progress is now happening!	
Rolf Morel Intel	MLIR now has a new <u>Area Team</u> to help guide its evolution, along with a <u>new</u> <u>organizational structure and charter</u> and <u>governance group</u> . The charter defines separate area groups: MLIR Core (the domain-independent infrastructure), and the dialects (like the machine learning-specific pieces). I am extremely thankful to everyone who is spending time to improve MLIR and work through these issues— such work has a profound impact on everyone building into the ecosystem as well as the downstream users.	

\*In bold: MLIR Area Team, a distinct governance effort



# Tensor Compiler Design Group – Modus Operandi

- Regular meetings
  - Bi-weekly group meetings
    - Agenda & notes: MLIR Tensor Compiler Design Group Meetings
    - ... only had three (3) so far
  - First Open Design Meeting scheduled late April (date TBC)



# Tensor Compiler Design Group – Modus Operandi

- Regular meetings
  - Bi-weekly group meetings
    - Agenda & notes: MLIR Tensor Compiler Design Group Meetings
    - ... only had three (3) so far
  - First Open Design Meeting scheduled late April (date TBC)
- In practice, members bring in-the-pipeline, community-relevant topics
  - In essence *pre-RFCs* workgroup provides space for quick iteration
    - Rapid top-of-mind responses for low overhead feedback
  - Next: posted as RFC, with same consensus process as any other RFC



# Tensor Compiler Design Group – Modus Operandi

- Regular meetings
  - Bi-weekly group meetings
    - Agenda & notes: MLIR Tensor Compiler Design Group Meetings
    - ... only had three (3) so far
  - First Open Design Meeting scheduled late April (date TBC)
- In practice, members bring in-the-pipeline, community-relevant topics
  - In essence *pre-RFCs* workgroup provides space for quick iteration
    - Rapid top-of-mind responses for low overhead feedback
  - Next: posted as RFC, with same consensus process as any other RFC
- All workgroup-generated documentation is public

MLIR Tensor Compiler Design Group - Overview document



# Tensor Compiler Design Group – Progress on Vector

MLIR Tensor Compiler Design Group - Vector Dialect: Refactoring + Re-design ideas

RFCs which benefitted from live discussion:

[RFC] Allow pointers as element type of `vector`

- Upshot: VectorElementTypeInterface with semantics of only allowing "atomic" elements

[RFC] Improving gather codegen for Vector Dialect

- Addresses abstraction gap which lead to early loss of structured indexing

[RFC] Generalize tiling to operate on ShapedType

- Proposes extending infrastructure for tiling/blocking beyond linalg-on-tensor/memref
- In-the-pipeline RFC on reducing references to LLVM LangRef in vector dialect docs

### Tensor Compiler – other recent significant changes

Transpose attribute for Linalg matmul operations

- Linalg.matmul (and batch\_matmul) now have an affine\_maps attribute

 $\square Introduce linalg.contract: D[J] = (\bigoplus_{(|A \cup |B|) \setminus J)} A[|A|]*B[|B|]) \oplus C[J]$ 

- Op generalizing all contraction ops (e.g. matmul variants and matvec and dot and ...)

Extend Linalg elemwise named ops semantics

– linalg.elementwise with affine\_maps replacing linalg.elem\_wise\_{unary,binary}

Move `tensor.pack` and `tensor.unpack` into Linalg

– Fix layering issue; facilitate interactions with linalg ops; allow for packing on memrefs



### Topics we are looking to make progress on

- 1. **[**RFC] Should we restrict the usage of 0-D vectors in the Vector dialect?
- 2. Remove ops with overlapping functionality from Vector, e.g. extractelement -> extract
- 3. Various issues in Vector regarding consistent naming and semantics
- 4. Pros & cons of splitting vector dialect into high-level and low-level parts
- 5. Vector vs Tensor types another go at clarifying their distinctive roles



# Topics we are looking to make progress on

- 1. <a>[RFC] Should we restrict the usage of O-D vectors in the Vector dialect?</a>
- 2. Remove ops with overlapping functionality from Vector, e.g. extractelement -> extract
- 3. Various issues in Vector regarding consistent naming and semantics
- 4. Pros & cons of splitting vector dialect into high-level and low-level parts
- 5. Vector vs Tensor types another go at clarifying their distinctive roles
- 6. Revisit how to attach layouts to tensors & vectors
- 7. Enshrine in the charter the significance of projected permutation affine\_maps
- 8. Pick up stalled progress on common-ing up linalg conv ops
- 9. Compute type on linalg ops given correspondence with imperfect loop nest
- 10. Op aliasing in linalg, e.g. linalg.matmul vs linalg.contract vs linalg.generic
  - Equiv. class perspective on matching w.r.t. **[**<u>[RFC]</u> <u>Linalg operation tree</u>

### Tensor Compiler Design Group – next steps

- More focus on documentation updates
  - Many dialect-level documents are outdated, as are the rationale ones
  - We would love help with this!



### Tensor Compiler Design Group – next steps

- More focus on documentation updates
  - Many dialect-level documents are outdated, as are the rationale ones
  - We would love help with this!
- Start formulating the Tensor Compiler-spanning, cross-dialect charter
  - New documents laying out intended coherent usage across dialects



### Tensor Compiler Design Group – next steps

- More focus on documentation updates
  - Many dialect-level documents are outdated, as are the rationale ones
  - We would love help with this!
- Start formulating the Tensor Compiler-spanning, cross-dialect charter
  - New documents laying out intended coherent usage across dialects
- Start codifying flows through the Tensor Compiler
  - E.g., integration tests spanning just the Tensor Compiler dialects
  - In addition to explicitly described flows in the charter



# Tensor Compiler Design Group – get in touch!



Alex Zinenko (<u>@ftynse</u>)



Quinn Dawkins (<u>@qed</u>)



Rolf Morel (@rolfmorel)





Renato Golin (<u>@rengolin</u>)



Javed Absar (<u>@javedabsar</u>)



Andrzej Warzynski (<u>@banach-space</u>)



Jacques Pienaar (<u>@jpienaar</u>)



Jakub Kuderski (@kuhar)



Diego Caballero (@dcaballe)



Matthias Springer (@matthias-springer)



Suraj Sudhir (<u>@sjarus</u>)



Kunwar Grover (@groverkss)

Tensor Compiler Design Group