

† Corresponding author

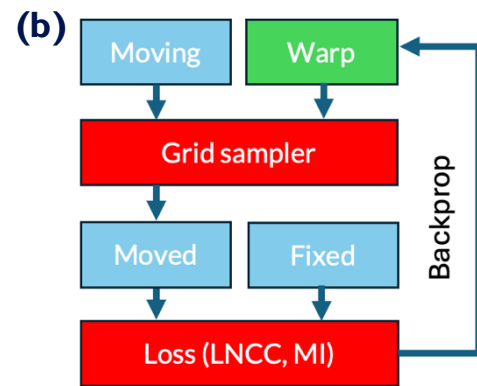
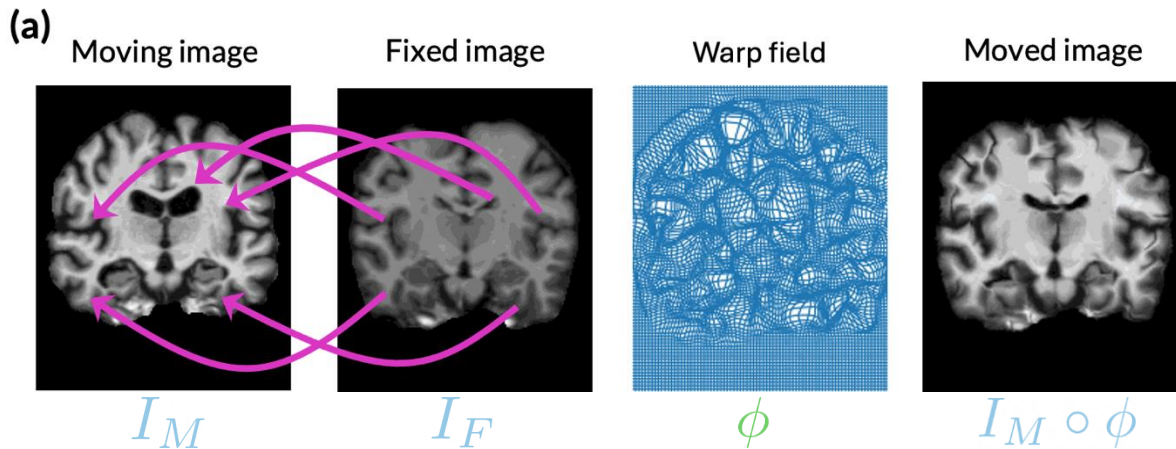
*These authors contributed equally

A Scalable Distributed Framework for Multimodal GigaVoxel Image Registration

Rohit Jena^{†,*}, Vedant Zope^{*}, Pratik Chaudhari, James Gee
{rjena, zope, pratikac}@seas.upenn.edu, gee@upenn.edu



Background



Process of establishing dense **Corresponding Matching** Between Fixed Image and Moving Image

$$\phi^* = \arg \min_{\phi \in \mathcal{G}} C(I_F, I_M \circ \phi) + \mathcal{R}(\phi)$$

Memory Bottlenecks in registration

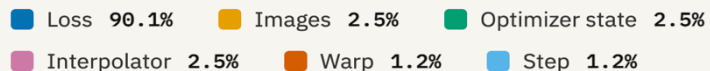
CROSS CORRELATION

Breakdown of GPU memory allocation per registration component

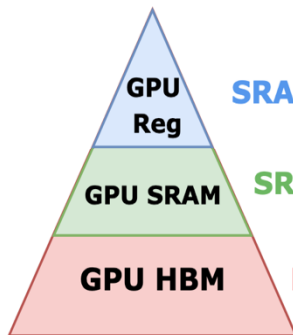
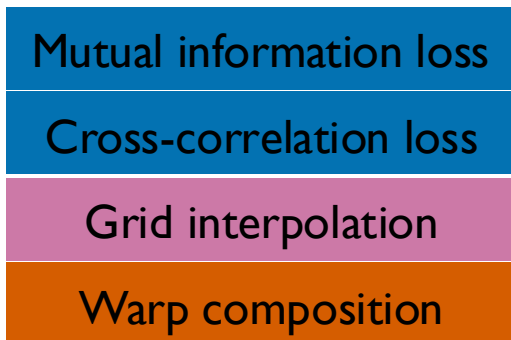


MUTUAL INFORMATION

Breakdown of GPU memory allocation per registration component



Intermediate tensors dominate runtime and memory



SRAM: 100TB/s

```
scalar_t val = x[index];  
val = cosf(cosf(val));
```

SRAM: 20TB/s

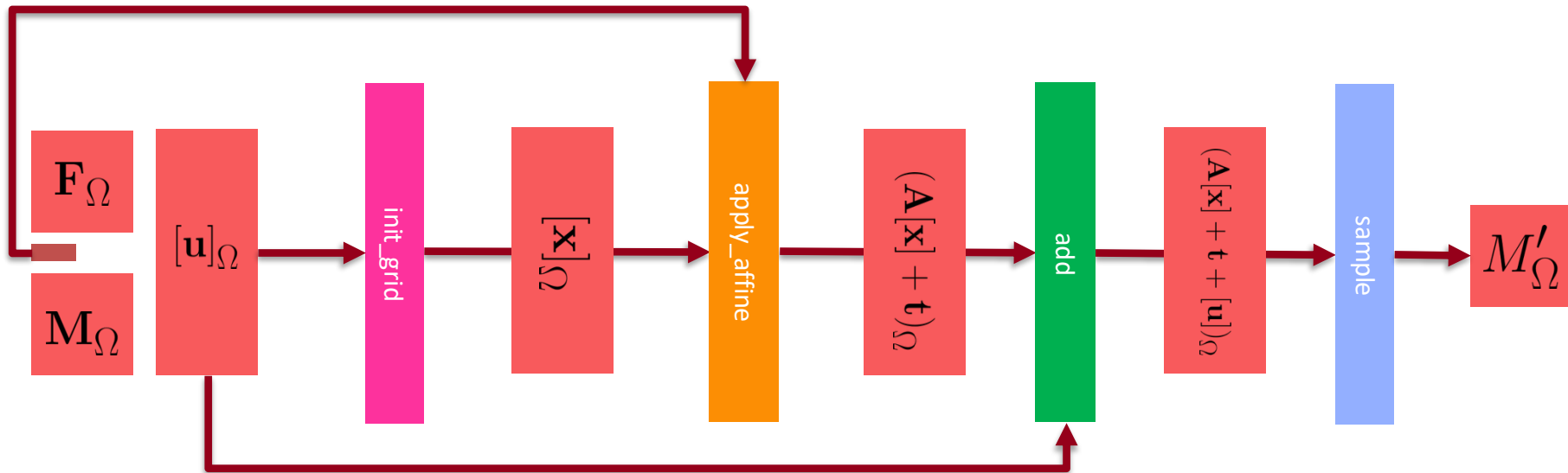
```
__shared__ scalar_t cu_sum[THREADS];  
cu_sum[threadIdx.x] += cu_sum[threadIdx.x * 2];
```

HBM: 1TB/s

```
x = torch.randn(5, 5).cuda();
```

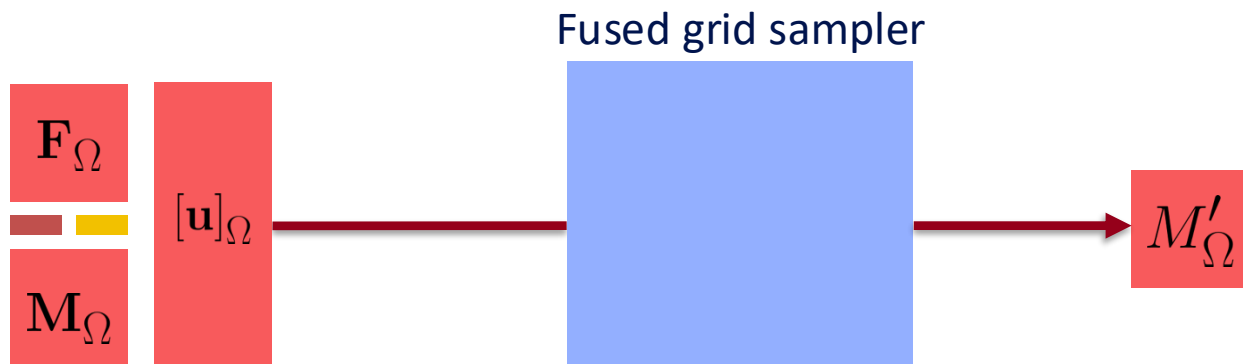
Grid Sampler/Interpolator

Composite Grid Sampler



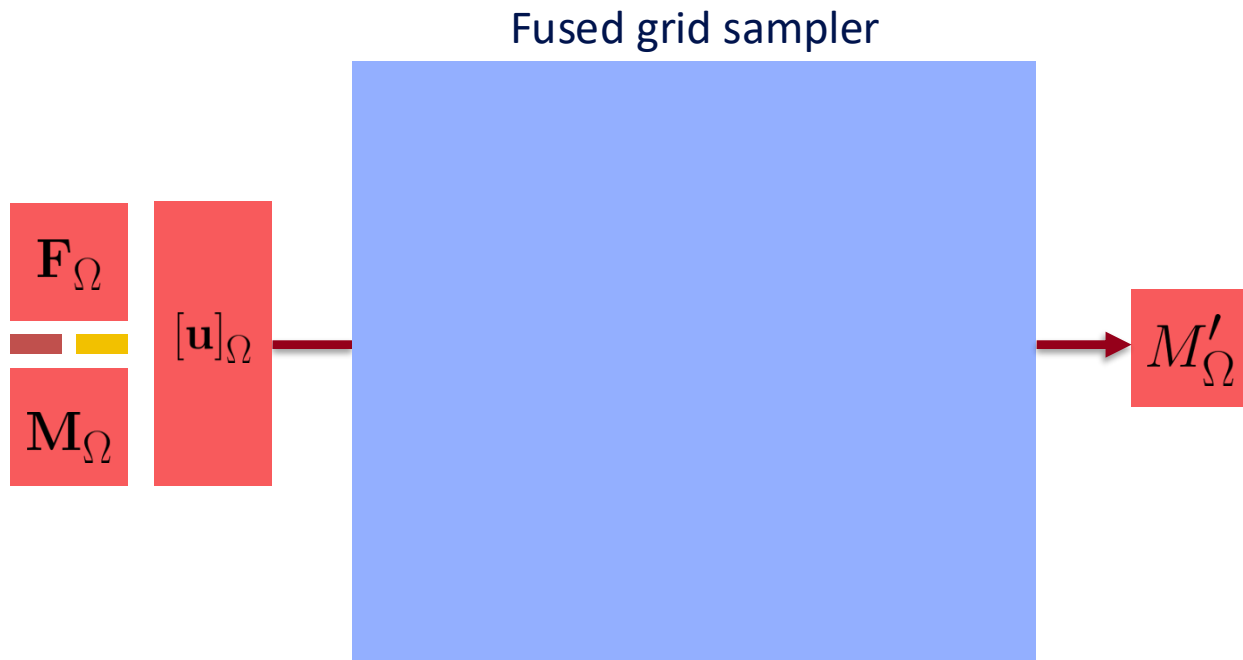
$$\text{grid_sampler}(I; A, t, [u])(x) = \text{sample}(I, \text{apply_affine}(A, t, \text{init_grid}(x)) + u)$$

Composite Implicit Grid Sampler



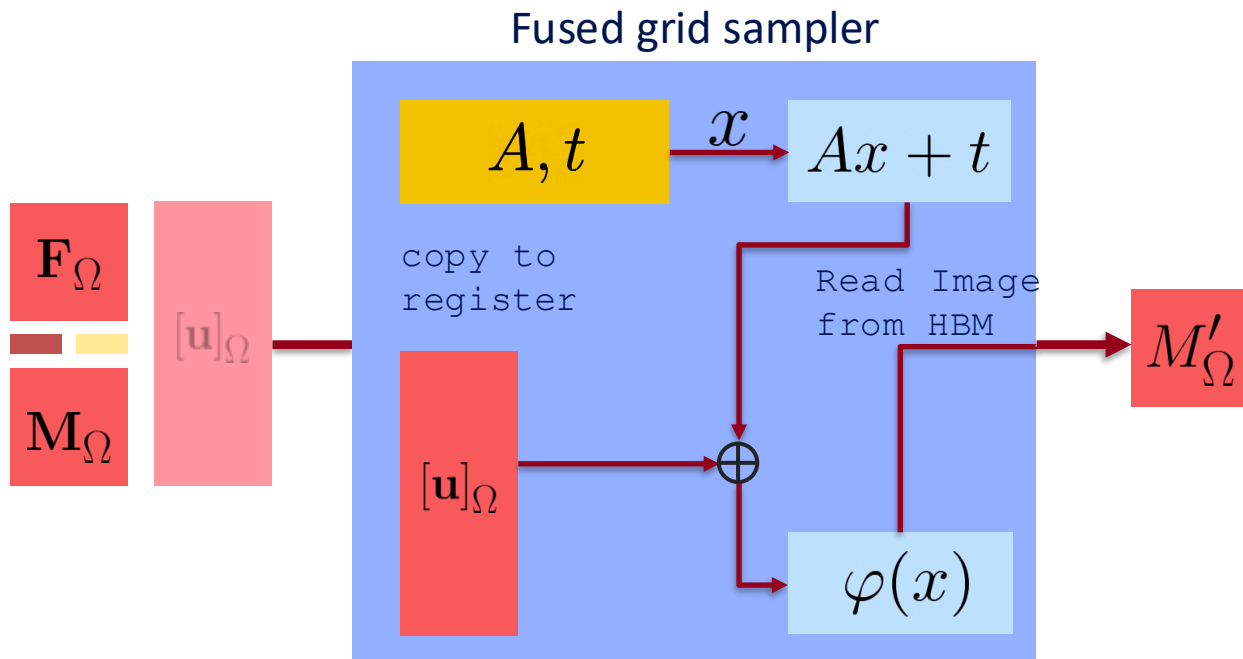
$$\text{fused_grid_sampler}(I; A, t, [u], S, x_{\text{bounds}})(x) = I(Ax + t + Su(x))$$

Composite Implicit Grid Sampler



$$\text{fused_grid_sampler}(I; A, t, [\mathbf{u}], S, x_{\text{bounds}})(x) = I(Ax + t + Su(x))$$

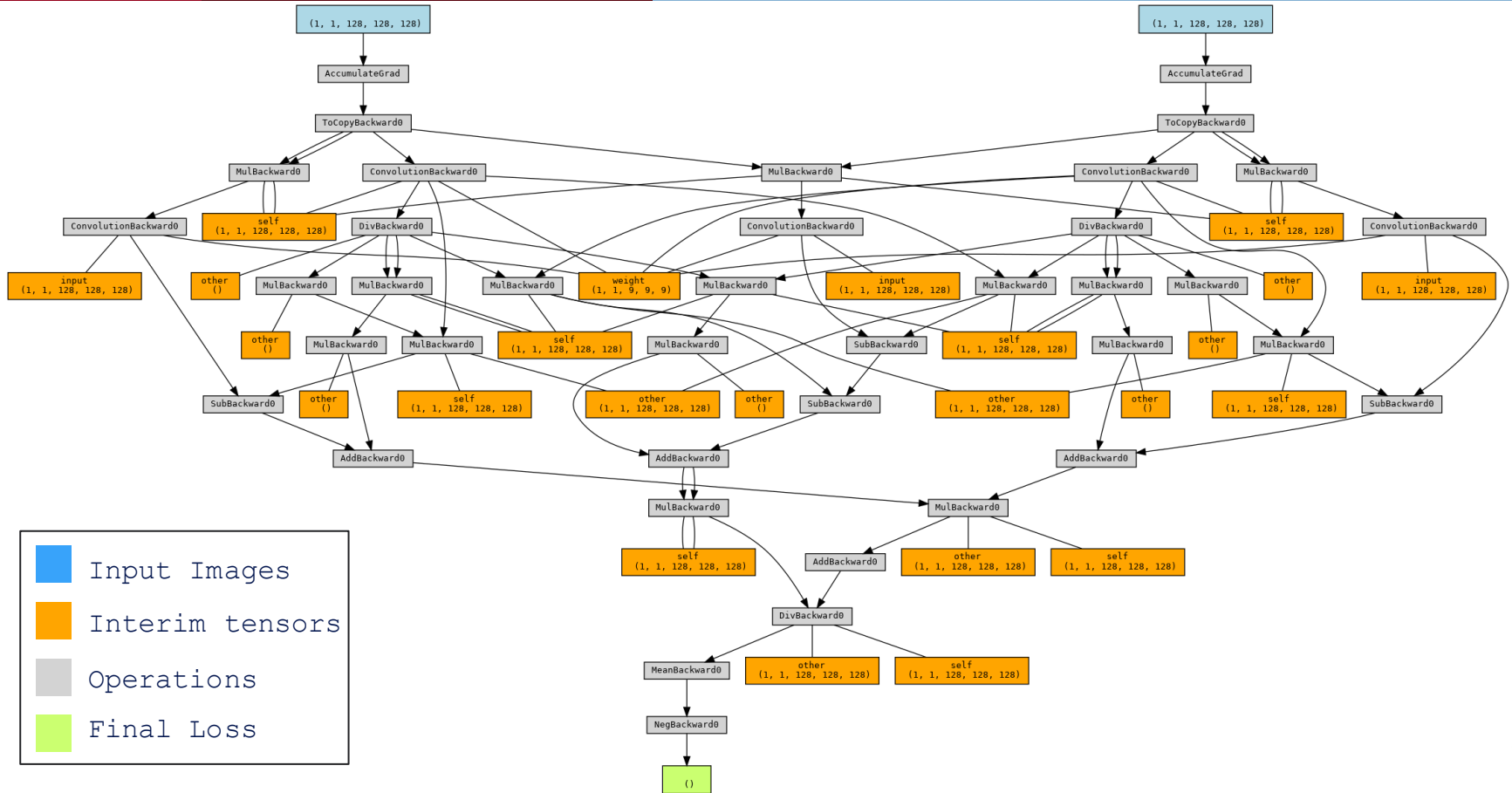
Composite Implicit Grid Sampler



$$\text{fused_grid_sampler}(I; A, t, [u], S, x_{\text{bounds}})(x) = I(Ax + t + Su(x))$$

Similarity Functions

Computational graph of the Vanilla LNCC



Cross Correlation Loss

 F_Ω M_Ω

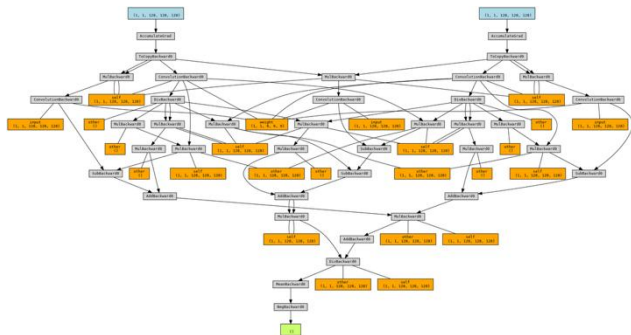
Vanilla LNCC

State Tensors + Intermediates

$\{F, M, F^2, M^2, FM,$
 $\mu_F, \mu_M, \mu_F^2, \mu_M^2 \dots\}$

Stores:

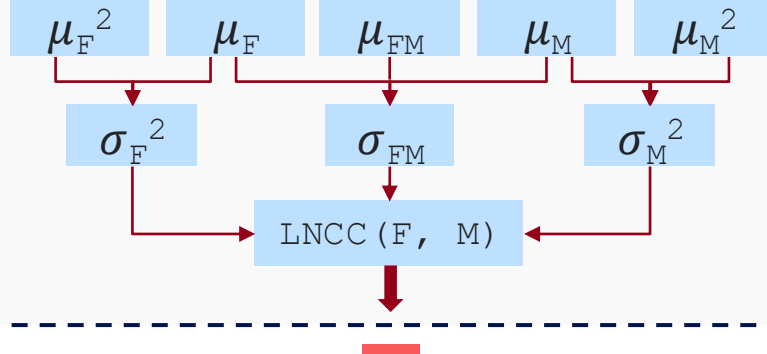
- $\sim 16N$ intermediates (forward)
- $\sim 16N$ additional memory (backward)

 F_Ω M_Ω

Fused LNCC

Analytical reformulation + kernel fusion

Stores only convolved states



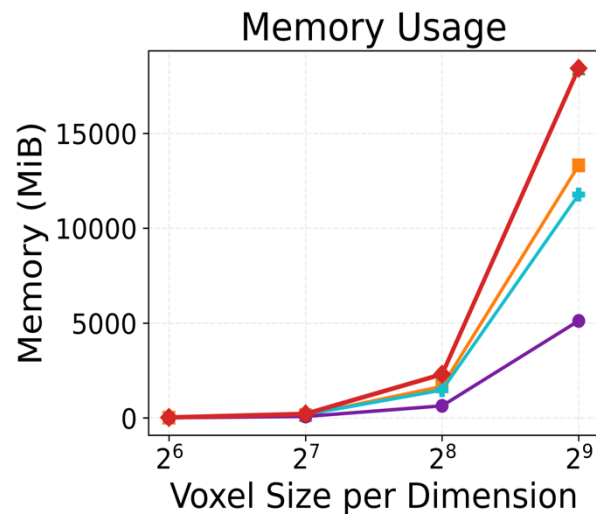
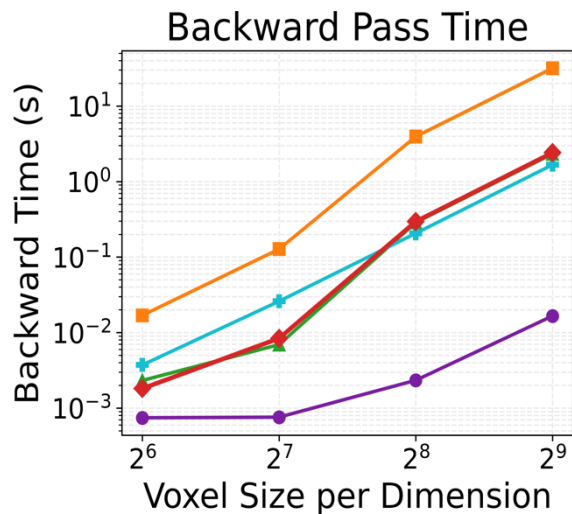
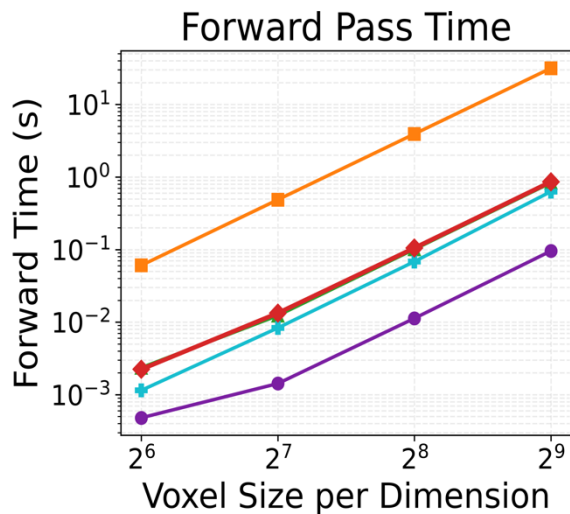
Stores:

- $\sim 5N$ intermediates (forward)
- No additional memory (backward)

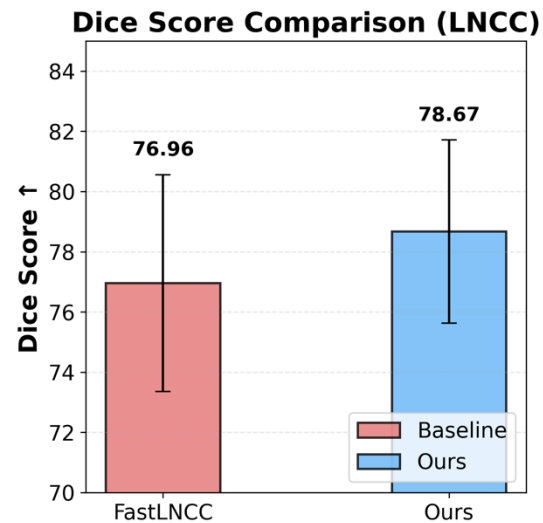
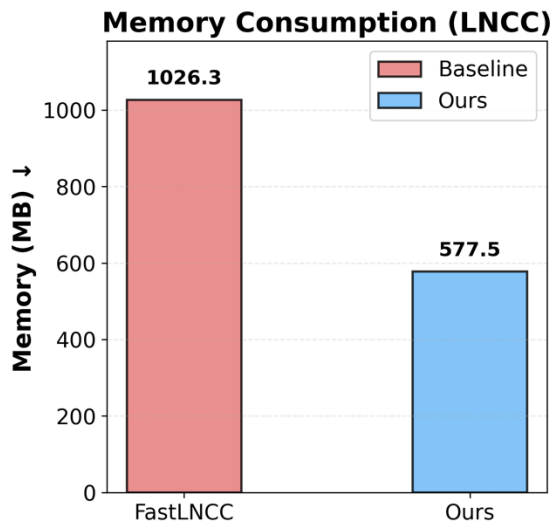
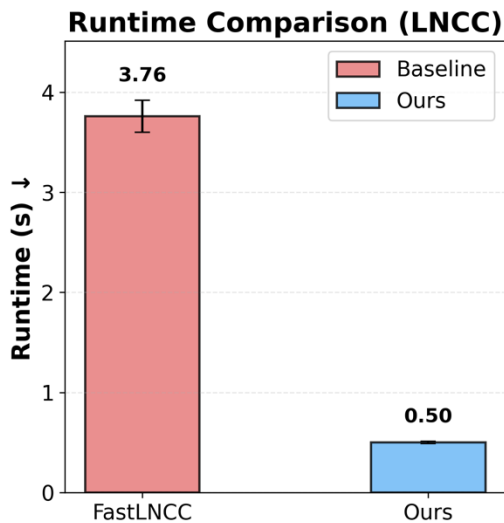
Fused Kernels Speedup (CC)

Ablation on CC Backends

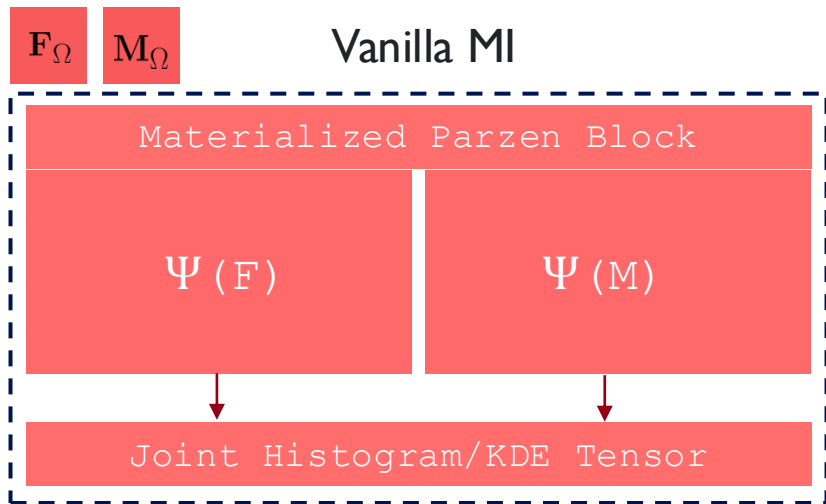
FireANTs torch.compile VoxelMorph Fast LNCC Ours



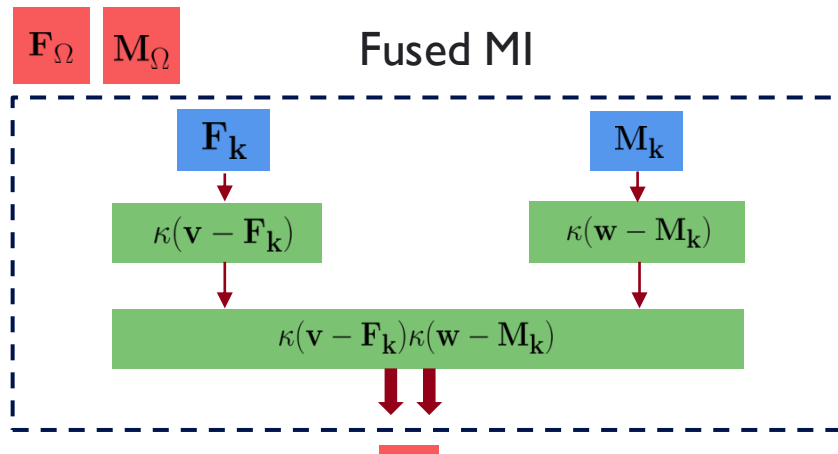
Iterative optimization becomes real-time



Mutual Information Loss



- Store full histogram in HBM
- Backprop requires storing all intermediates
- $O(N)$ HBM overhead

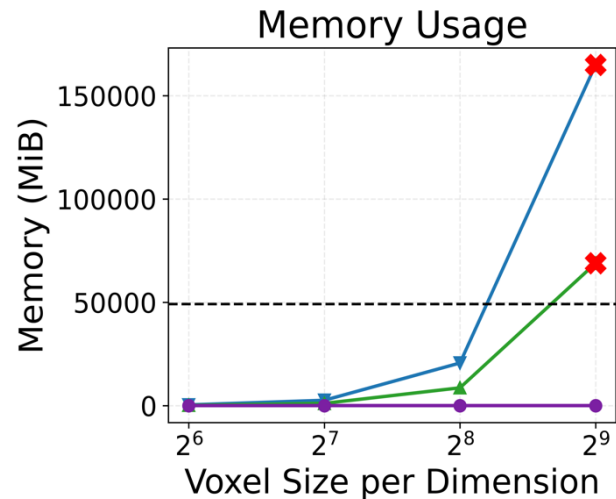
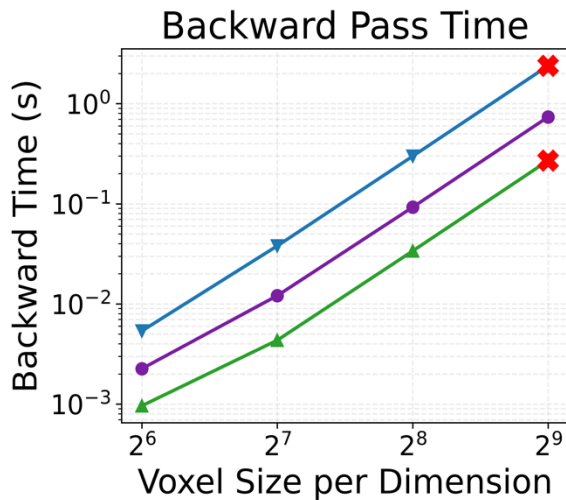
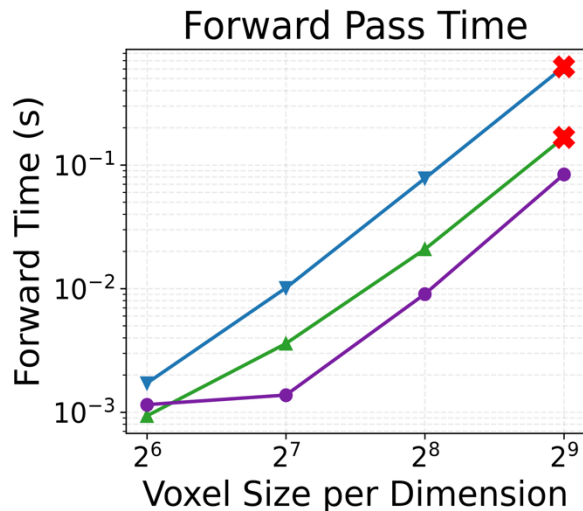


- Kernel contribution are computed per-thread and accumulated in shared memory
- Backward computed on-the-fly
- $O(1)$ HBM overhead (holding B constant)

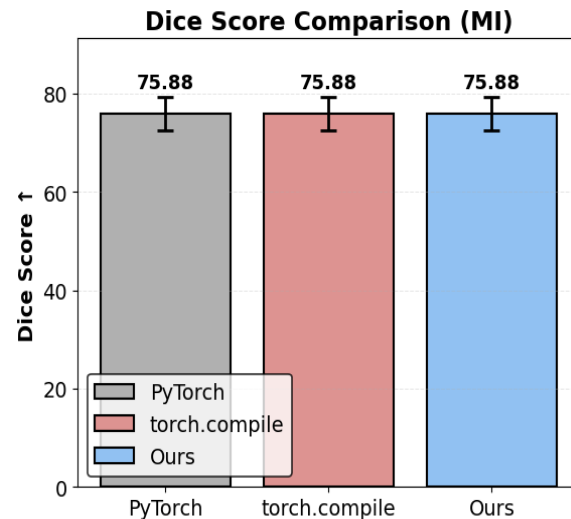
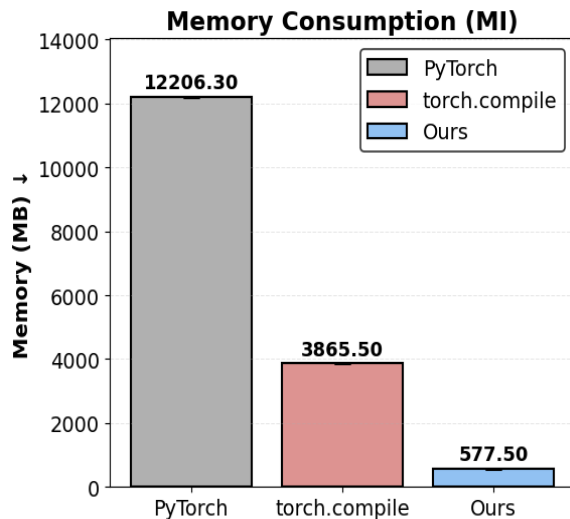
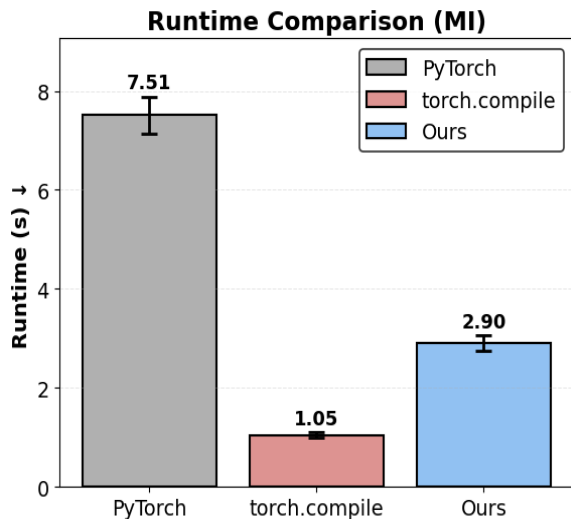
Fused Kernels Speedup (MI)

Ablation on Mutual Information Backends

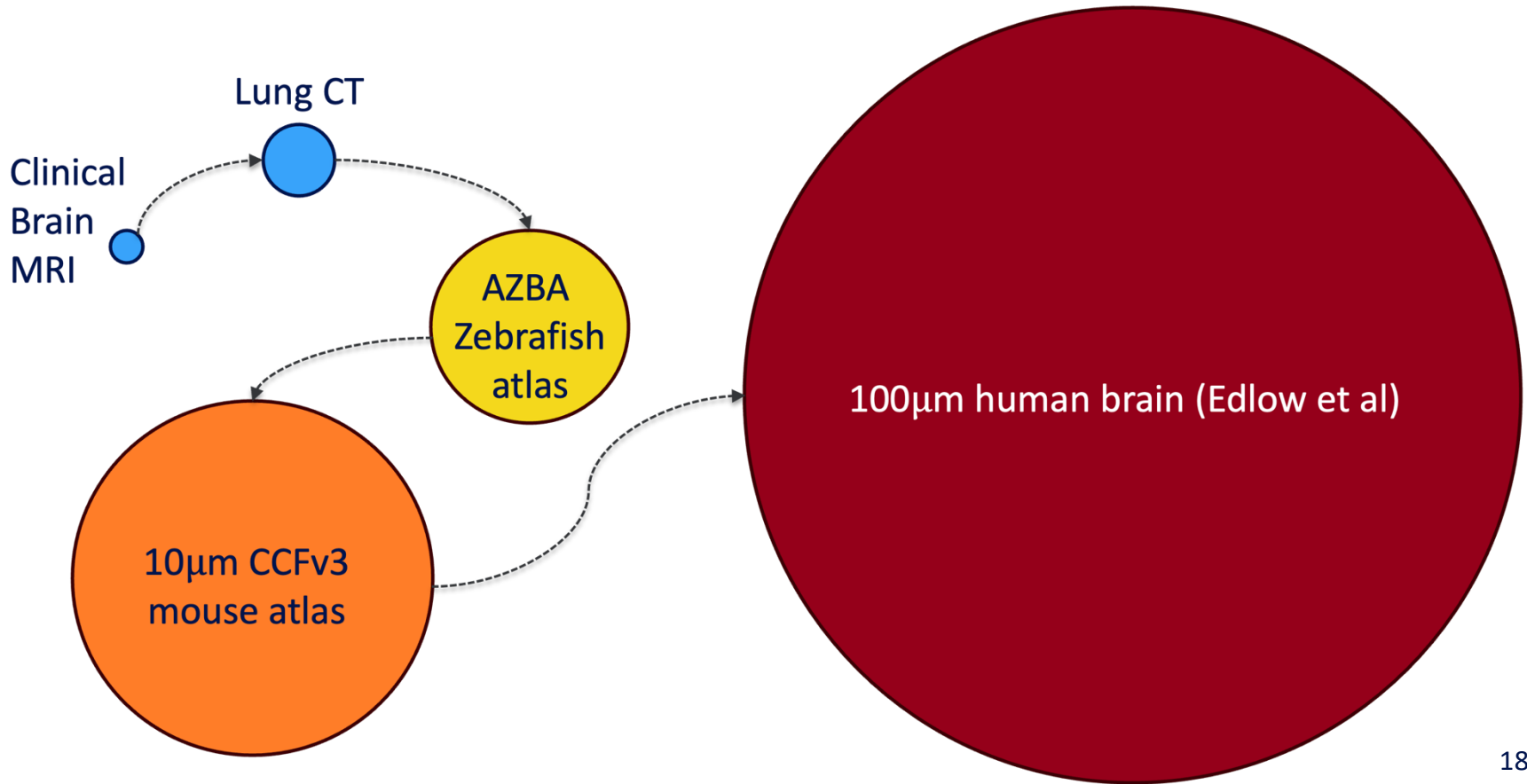
— PyTorch — torch.compile — Ours



Iterative optimization becomes real-time

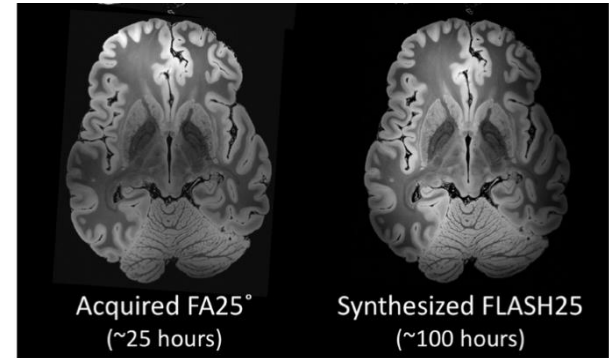
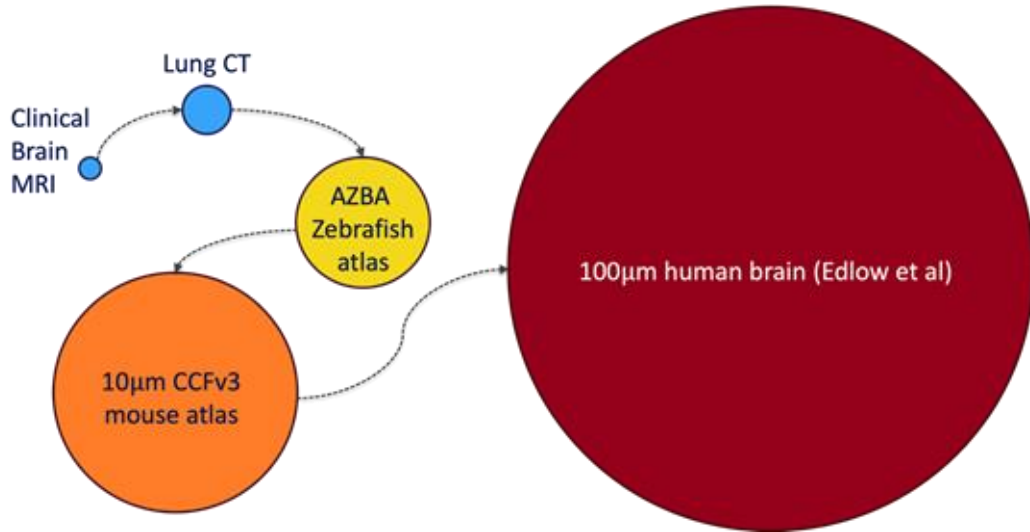


Imaging capabilities have scaled rapidly



Imaging capabilities have scaled rapidly

What happens when Image is too large?

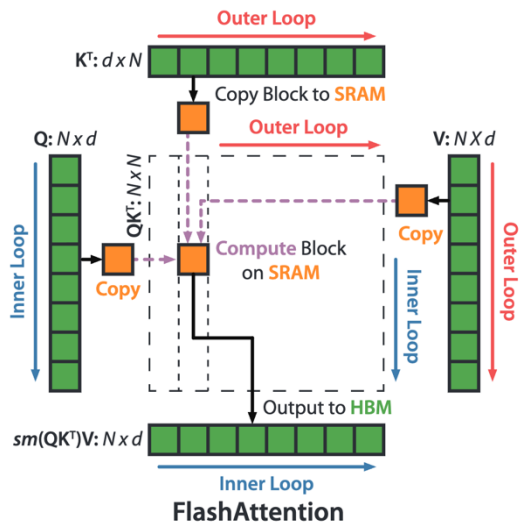


$[u]_{\Omega} \approx 44\text{GB VRAM}$

Extending image registration to multiple GPUs

Background: Existing work

Two key ideas: Tiling, Kernel Fusion, Tensor Parallelism



Megatron-LM and Megatron Core

GPU-optimized library for training transformer models at scale

docs latest release 0.15.0 license Apache

Background: Existing work

Two key ideas: **Tiling**, **Tensor Parallelism**

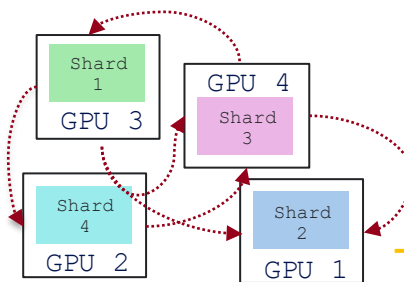


Why existing approaches don't scale?



Random access:

Tiling ❌

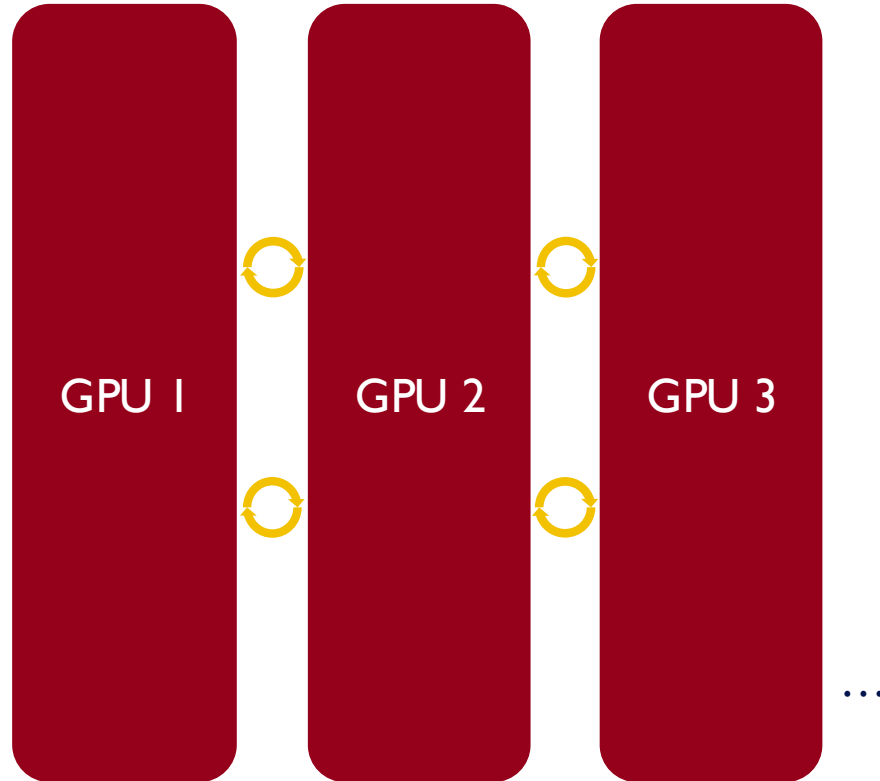
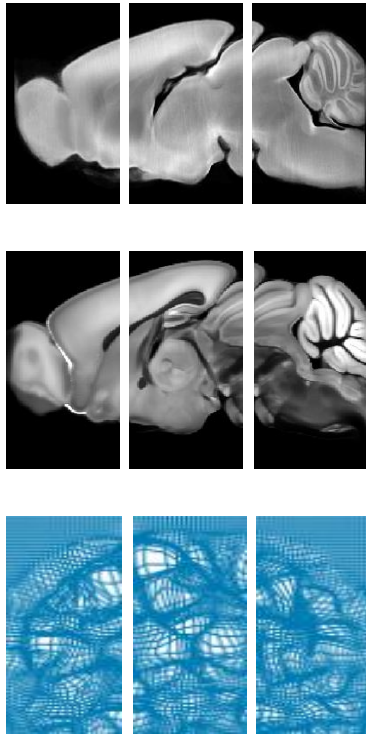


Random access
+
boundary sync:

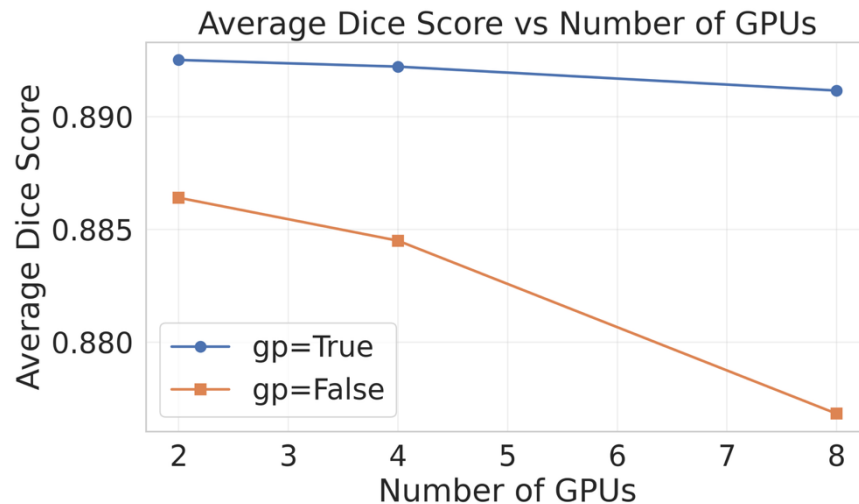
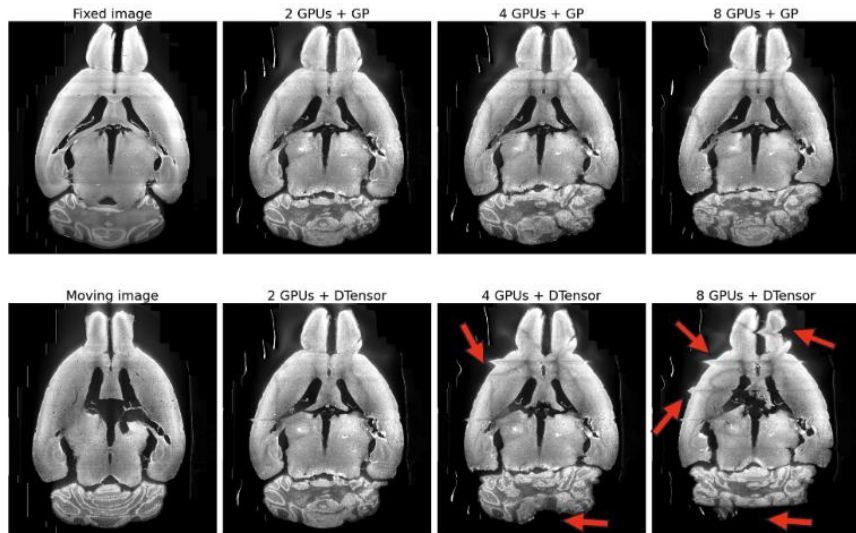
Tensor Parallelism ❌

Grid Parallel

Grid Parallel



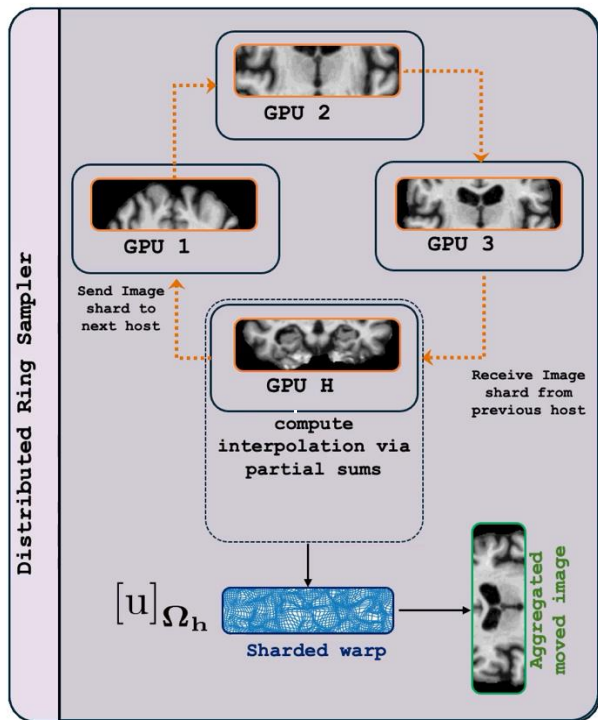
Why Grid Parallel?



Distributed Ring Sampler

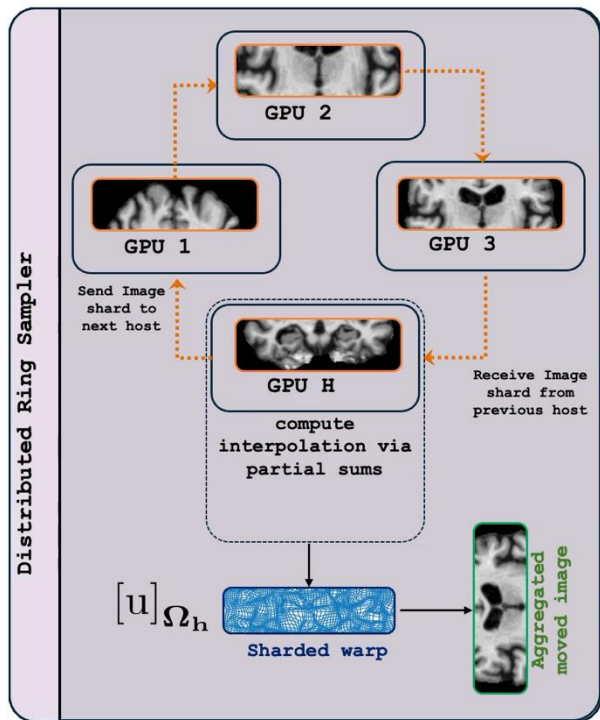
Distributed Ring Sampler

Ring Communication
(All-to-All via sequential passing)

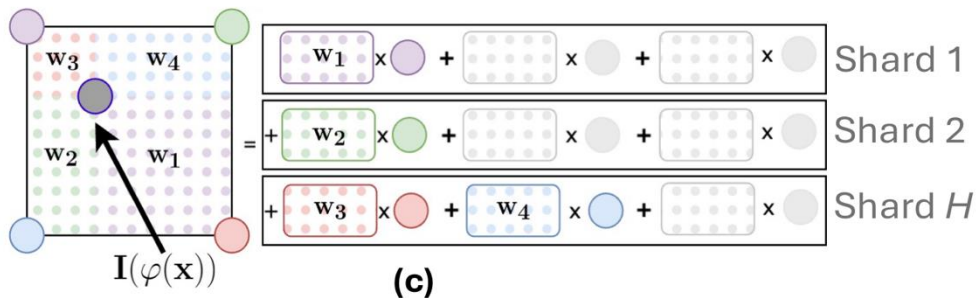


Distributed Ring Sampler

Ring Communication
(All-to-All via sequential passing)

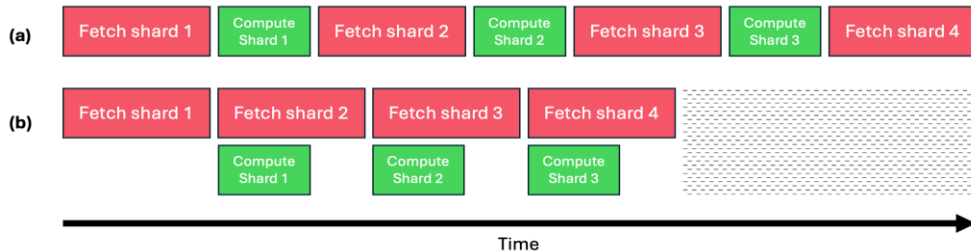


Interpolation requires cross-shard data dependencies



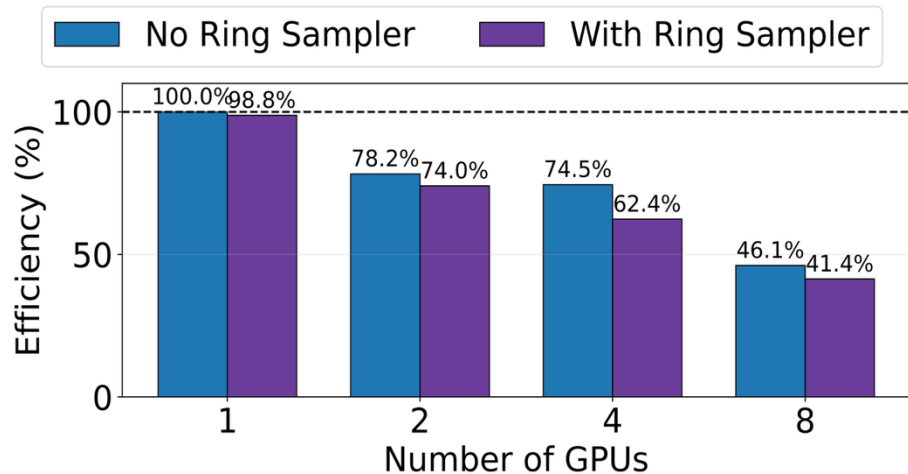
Partial sums \rightarrow final interpolation

Each shard computes weighted contributions (w_i)

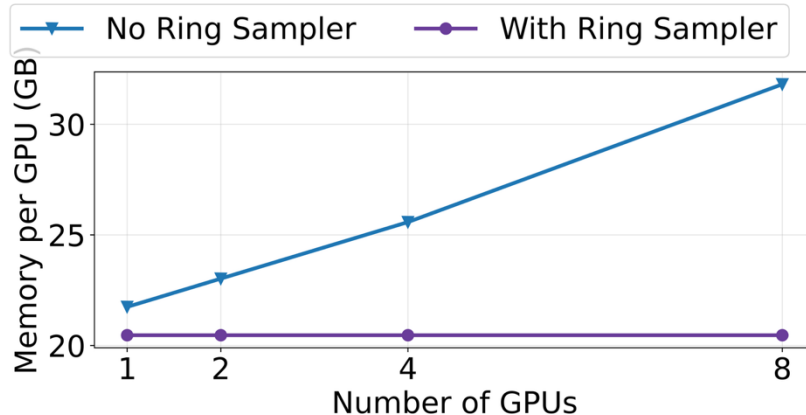


Scaling Results

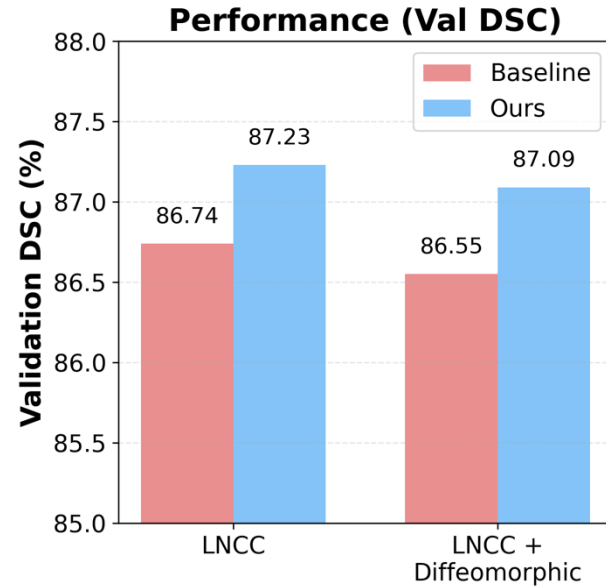
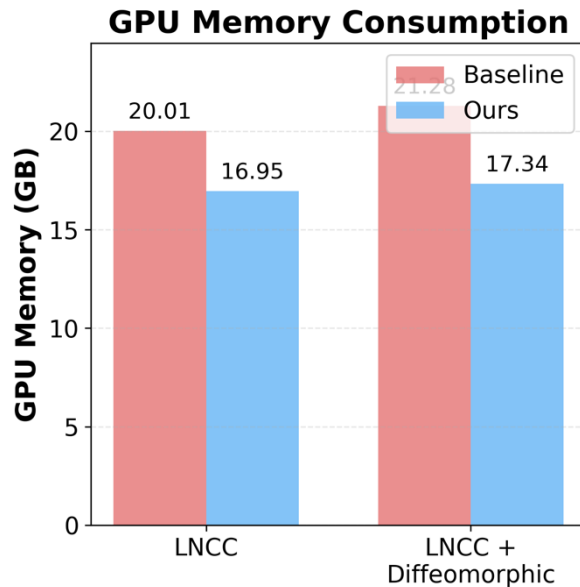
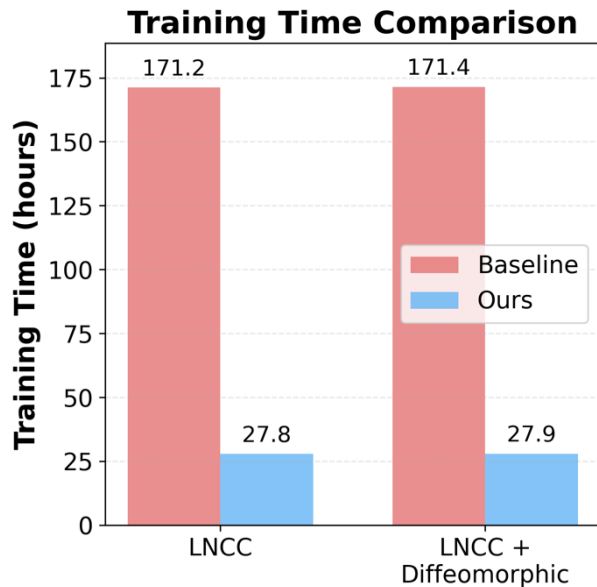
Weak Scaling



Memory Consumption per GPU



Deep Learning training speeds up by 6x



TransMorph

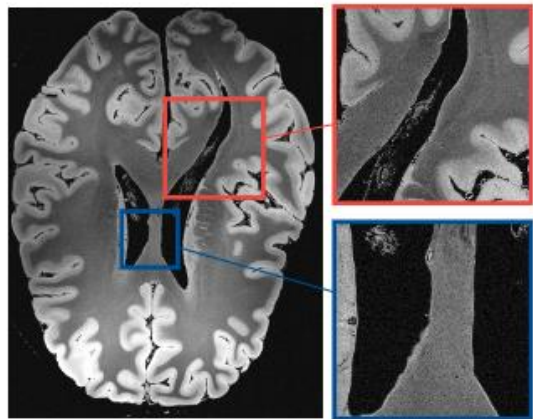
Multimodal Gigavoxel registration

Ex vivo (Edlow et al.) : **100** μm FLASH imaging
1760 \times **1760** \times **1278** \approx **11.8B** voxels

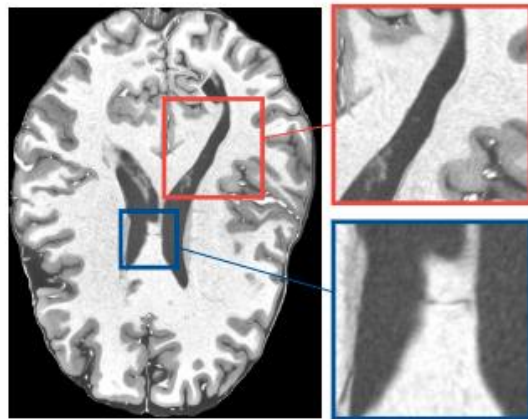
In vivo (Lusebrink et al.) : **250** μm T1-weighted
800 \times **800** \times **800**

Using **8** A6000 GPUs in **\sim 1 minute**

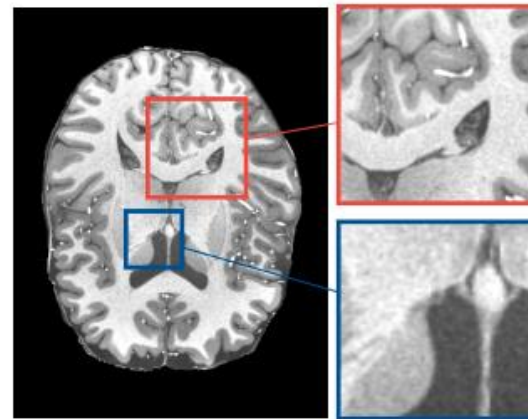
Qualitative Results



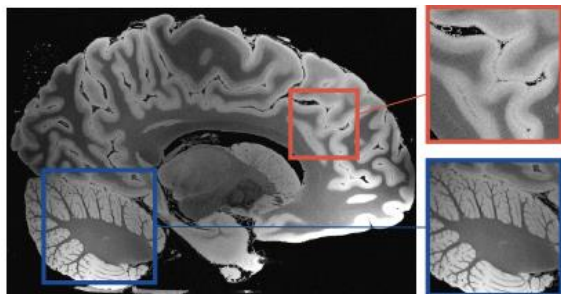
Fixed ($100\mu\text{m}$ FLASH)



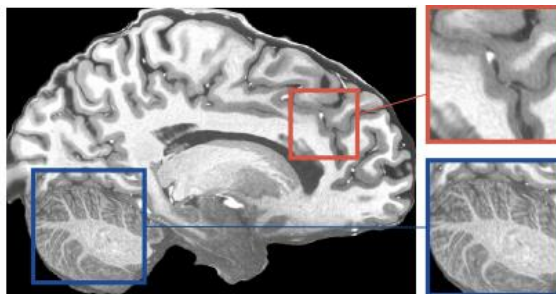
Moved ($250\mu\text{m} \rightarrow 100\mu\text{m}$)



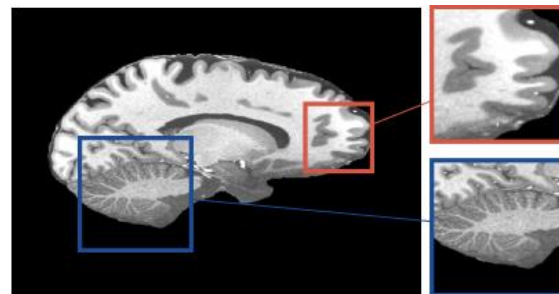
Moving ($250\mu\text{m}$ T1)



Fixed ($100\mu\text{m}$ FLASH)



Moved ($250\mu\text{m} \rightarrow 100\mu\text{m}$)



Moving ($250\mu\text{m}$ T1)

Thank you!

Q&A

Welcome to join me for a discussion during the poster session!

 April 25, 2026,

 2:15 PM – 4:45 PM (EDT)

 Poster Session 6