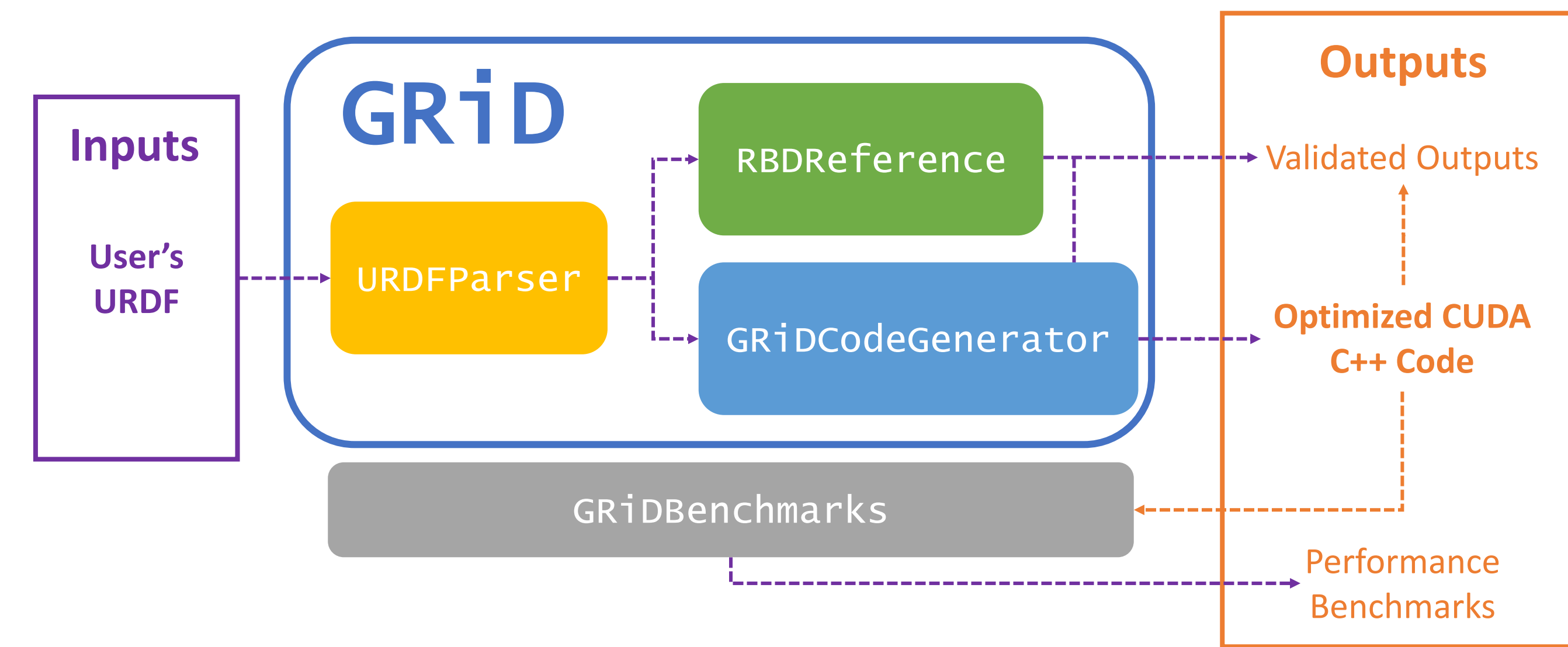


GRiD: GPU-Accelerated Rigid Body Dynamics with Analytical Gradients

Brian Plancher¹, Sabrina M. Neuman¹, Radhika Ghosal¹, Scott Kuindersma^{1,2}, Vijay Janapa Reddi¹

1: Harvard University John A. Paulson School of Engineering and Applied Sciences, 2: Boston Dynamics

The Big Picture:



github.com/robot-acceleration/GRiD

GRiD is a header-only, modular, open-source, GPU-accelerated library for rigid body dynamics with analytical gradients. Key features include:

- **URDF parsing & code generation** to deliver optimized dynamics kernels that expose GPU-friendly computational patterns
 - E.g., Leverages both fine-grained parallelism within each computation & coarse-grained parallelism between computations
- Delivers **end-to-end computational speedups** through algorithmic refactoring
- **Modular, open-source, and header-only**

GRiD currently supports:

- Prismatic, fixed, and revolute joints
- **ID, FD, M^{-1}**
- **∇ ID, ∇ FD** with respect to q, \dot{q}, u

This material is based upon work supported by the National Science Foundation (under Grant DGE1745303 and Grant 2030859). Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors and may not reflect those of the funding organizations.

Design Optimizations:

We **re-factored** algorithms to better leverage the GPU's strengths by:

- Exposing more **natural parallelism** (especially across branches)
- Reducing work done in serial loops
- Leveraging topology driven **sparsity patterns** in matrices
- Unifying computational operations

Algorithm 1 ∇ RNEA-F($\dot{q}, v, a, f, X, S, I$) $\rightarrow \partial c / \partial u$

```

1: for frame  $i = 1 : N$  do
2:    $\frac{\partial v_i}{\partial u} = {}^i X_{\lambda_i} \frac{\partial v_{\lambda_i}}{\partial u} + \begin{cases} ({}^i X_{\lambda_i} v_{\lambda_i}) \times S_i & u \equiv q \\ S_i & u \equiv \dot{q} \end{cases}$ 
3:    $\frac{\partial a_i}{\partial u} = {}^i X_{\lambda_i} \frac{\partial a_{\lambda_i}}{\partial u} + \frac{\partial v_{\lambda_i}}{\partial u} \times S_i \dot{q}_i + \begin{cases} ({}^i X_{\lambda_i} a_{\lambda_i}) \times S_i \\ v_i \times S_i \end{cases}$ 
4:    $\frac{\partial f_i}{\partial u} = I_i \frac{\partial a_i}{\partial u} + \frac{\partial v_i}{\partial u} \times * I_i v_i + v_i \times * I_i \frac{\partial v_i}{\partial u}$ 

```

Algorithm 2 ∇ RNEA-F-GRiD($\dot{q}, v, a, f, X, S, I$) $\rightarrow \partial f / \partial u$

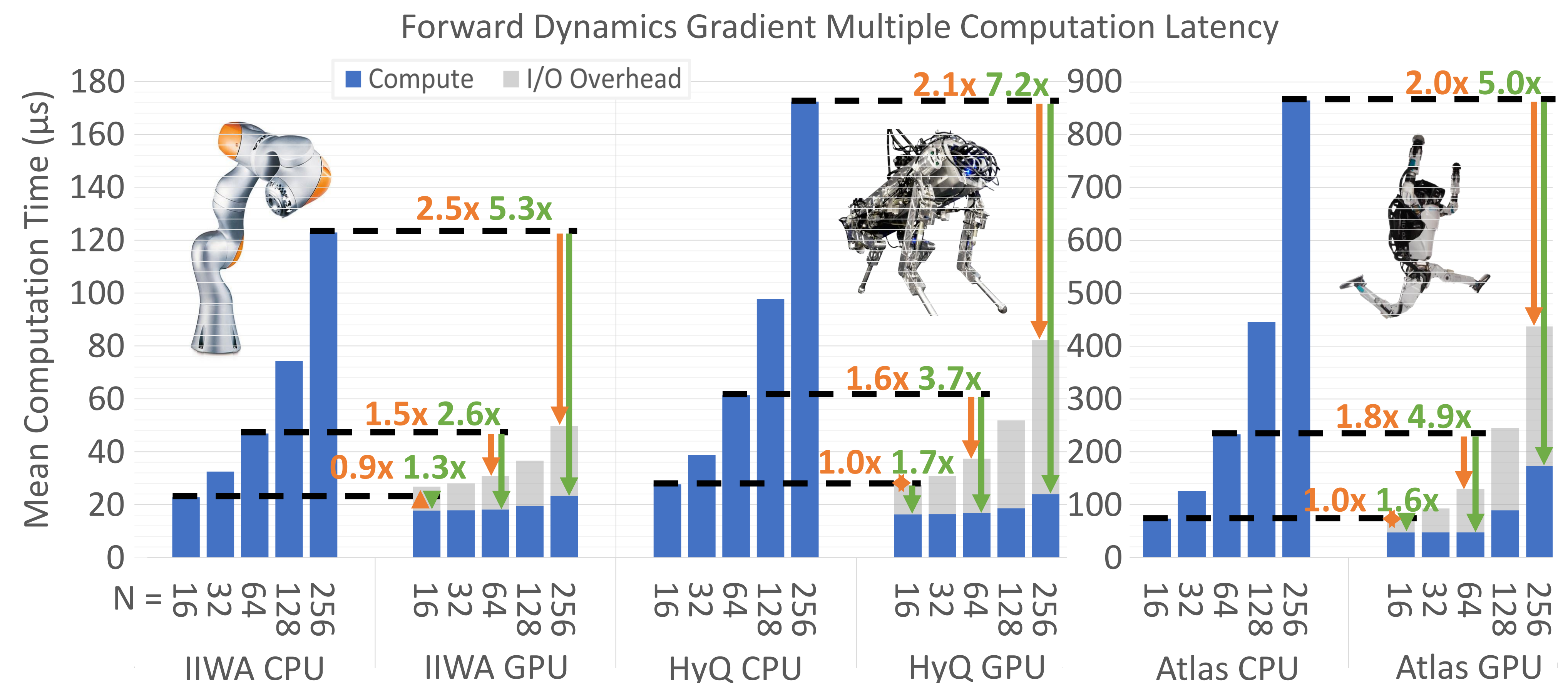
```

1: for frame  $i = 1 : n$  in parallel do
2:    $\alpha_i = {}^i X_{\lambda_i} v_{\lambda_i}$     $\beta_i = {}^i X_{\lambda_i} a_{\lambda_i}$     $\gamma_i = I_i v_i$ 
3:    $\alpha_i = \alpha_i \times S_i$     $\beta_i = \beta_i \times S_i$     $\delta_i = v_i \times S_i$ 
4: for level  $l = 0 : l_{max}$  do
5:   for frame  $i \in l$  in parallel do
6:      $\frac{\partial v_i}{\partial u} = {}^i X_{\lambda_i} \frac{\partial v_{\lambda_i}}{\partial u} + \begin{cases} \alpha_i & u \equiv q \\ S_i & u \equiv \dot{q} \end{cases}$ 
7: for frame  $i = 1 : n$  in parallel do
8:    $\rho_i = \frac{\partial v_{\lambda_i}}{\partial u} \times S_i \dot{q}_i + \begin{cases} \beta_i \\ \delta_i \end{cases}$ 
9: for level  $l = 0 : l_{max}$  do
10:  for frame  $i \in l$  in parallel do
11:     $\frac{\partial a_i}{\partial u} = {}^i X_{\lambda_i} \frac{\partial a_{\lambda_i}}{\partial u} + \rho_i$ 
12: for frame  $i = 1 : n$  in parallel do
13:    $\frac{\partial f_i}{\partial u} = \frac{\partial v_i}{\partial u} \times * \gamma_i$     $\eta_i = v_i \times * I_i$ 
14:    $\frac{\partial f_i}{\partial u} = \frac{\partial f_i}{\partial u} + I_i \frac{\partial a_i}{\partial u} + \eta_i \frac{\partial v_i}{\partial u}$ 

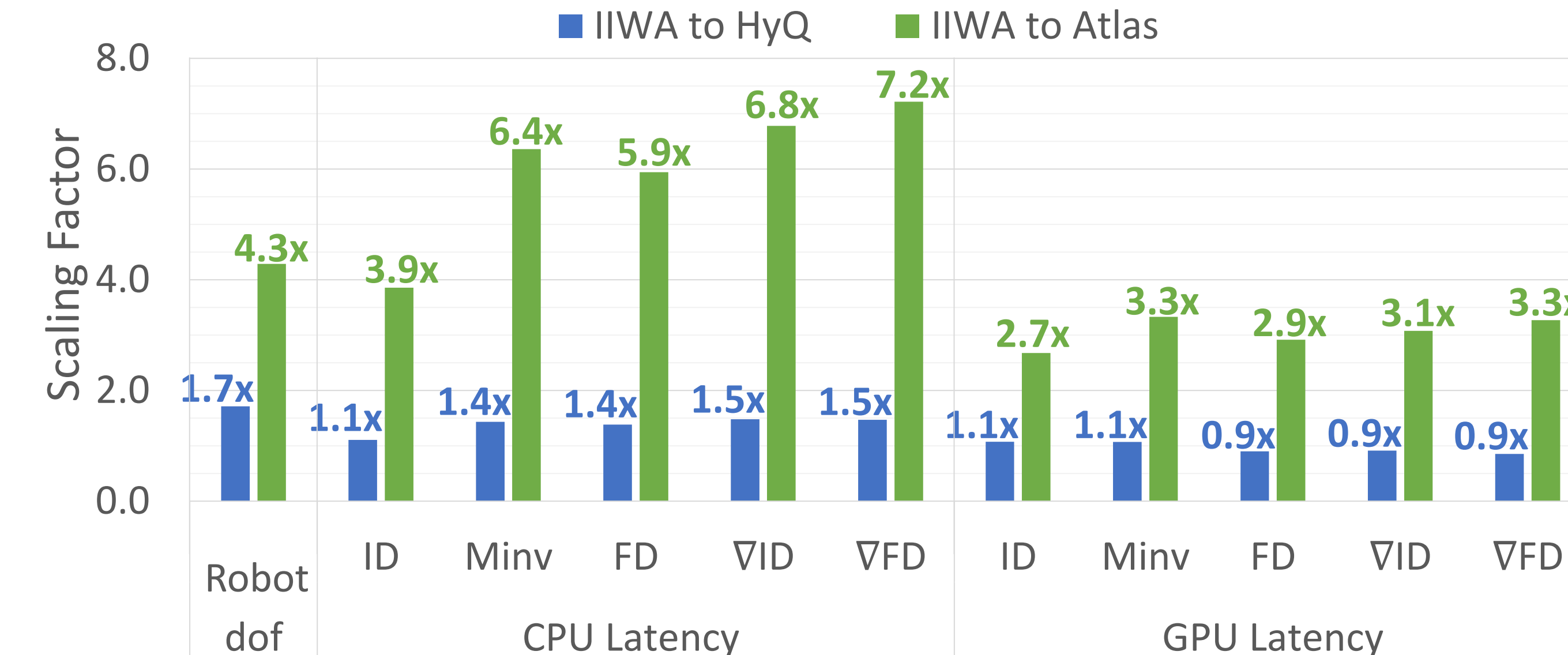
```

Performance Results:

- Benchmarked against Pinocchio, a state-of-the-art CPU library
 - Pinocchio supports optimized CPU code generation of rigid body dynamics & analytical gradients
- **GRiD scales well to complex robots and multiple computations**
 - As much as a 7.2x computational speedup over the CPU
 - As much as a 2.5x speedup when accounting for I/O overhead



Single Computation Latency Scaling across Robots and Algorithms



Single Computation Latency (μ s)
(ID = Inverse Dynamics, M^{-1} = Direct Inverse of the Mass Matrix, FD = Forward Dynamics, and ∇ indicates the gradient of that algorithm)

	Algorithm	IIWA	HyQ	Atlas
CPU	ID	0.3	0.3	1.1
	M^{-1}	0.5	0.8	3.4
	FD	0.9	1.2	5.3
	∇ ID	1.4	2.1	9.8
	∇ FD	2.9	4.3	20.9
GPU	ID	3.0	3.2	8.0
	M^{-1}	5.2	5.6	17.4
	FD	7.7	6.9	22.4
	∇ FD	12.9	11.0	42.1

We used a high-performance workstation with a 3.8GHz eight-core Intel Core i7-10700K CPU and a 1.44GHz NVIDIA GeForce RTX 3080 GPU running Ubuntu 20.04 and CUDA 11.4.4 We compare timing results across three robot models: the 7 degrees-of-freedom (dof) Kuka LBR IIWA-14 manipulator, the 12 dof HyQ quadruped, and the 30 dof Atlas humanoid. For single computation and multiple computation latency, we took the average of one million, and one hundred thousand trials, respectively.