

# Approaches to Nonlinear Programming on GPU Architectures

**Alexis Montoison<sup>1</sup>, Sungho Shin<sup>2</sup>, François Pacaud<sup>3</sup>,**  
**Michel Schanen<sup>2</sup>, and Mihai Anitescu<sup>2,4</sup>**

<sup>1</sup>Polytechnique Montréal and GERAD, Canada

<sup>2</sup>Mathematics and Computer Science Division, Argonne National Laboratory

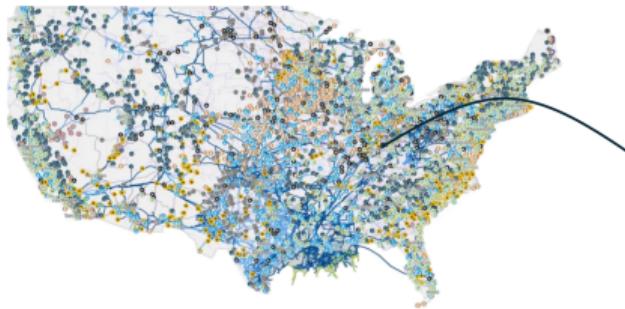
<sup>3</sup>Centre Automatique et Systèmes, Mines Paris - PSL

<sup>4</sup>Department of Statistics, University of Chicago

April 17, 2024  
GEEKS – Gurobi

# Takeaways

- Accelerated computing offers significant opportunities for systems engineering
  - A nonlinear solver on GPU is 10x faster than the state-of-the-art open-source solver on CPU [Shin, 2024]



Can optimize the operation of the **entire eastern interconnection** in <20 seconds

- A linear programming solver on GPU is comparable to state-of-the-art commercial solvers on CPU [Lu, 2024]
- We envision **expanding the application scope of optimization** with accelerated computing.
  - **Extremely large-scale** problems (coupled infrastructures, multi-stage, multiscale)
  - Optimization with **AI models** (e.g., surrogates) and **digital twins**

# Outline

1. Motivation
2. Nonlinear optimization framework
3. Implementing algebraic modeling systems for GPUs
4. Implementing nonlinear programming solvers for GPUs
  - 4.1. LiftedKKT
  - 4.2. HybridKKT
  - 4.3. Numerical results
5. Remaining challenges
6. Conclusion

# Outline

1. Motivation
2. Nonlinear optimization framework
3. Implementing algebraic modeling systems for GPUs
4. Implementing nonlinear programming solvers for GPUs
  - 4.1. LiftedKKT
  - 4.2. HybridKKT
  - 4.3. Numerical results
5. Remaining challenges
6. Conclusion

# Accelerated Computing in 2024

- ▶ Accelerated computing has **driven the success of AI** (e.g., GPT models have  $10^{12}$  pars).

## Accelerated Computing in 2024

- ▶ Accelerated computing has **driven the success of AI** (e.g., GPT models have  $10^{12}$  pars).
- ▶ Accelerated computing **empowers scientific computing** (e.g., fluid, climate, bioinformatics).

# Accelerated Computing in 2024

- ▶ Accelerated computing has **driven the success of AI** (e.g., GPT models have  $10^{12}$  pars).
- ▶ Accelerated computing **empowers scientific computing** (e.g., fluid, climate, bioinformatics).
- ▶ We're entering **exascale computing era** ( $10^{18}$  floating point operations per second).

Aurora Supercomputer @ Argonne



Mostly powered by GPUs

= 1 million ×

iPhone 14 Pro



# Accelerated Computing in 2024

- ▶ Accelerated computing has **driven the success of AI** (e.g., GPT models have  $10^{12}$  pars).
- ▶ Accelerated computing **empowers scientific computing** (e.g., fluid, climate, bioinformatics).
- ▶ We're entering **exascale computing era** ( $10^{18}$  floating point operations per second).

Aurora Supercomputer @ Argonne



Mostly powered by GPUs

= 1 million ×

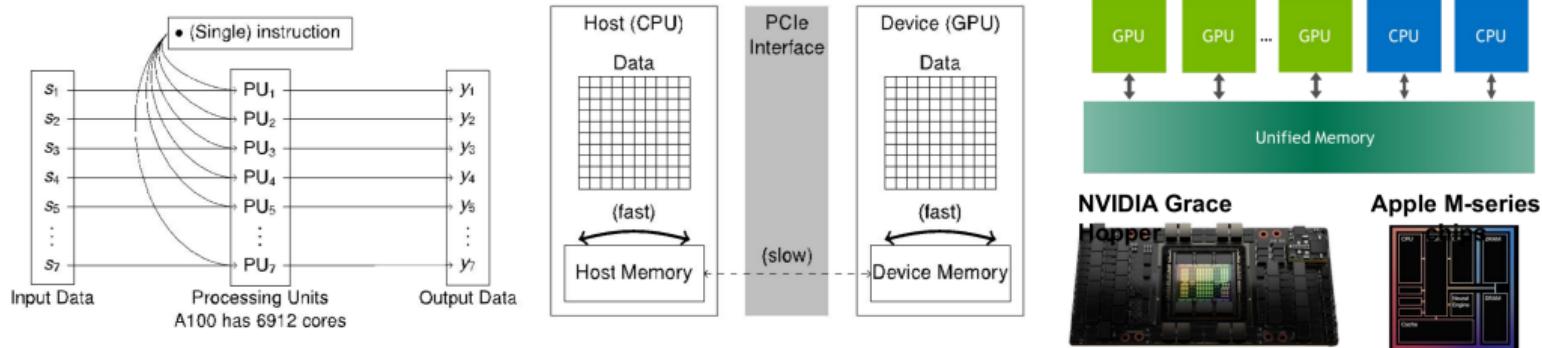
iPhone 14 Pro



Can we harness these capabilities in the realm of **classical nonlinear optimization** (e.g., energy infrastructures, optimal control, operations research)?

# How Do GPUs Work?

- Single Instruction, Multiple Data (**SIMD**) parallelism
- Dedicated device memory and slow interface: all data should reside in device memory only
- Emerging architectures employ **unified memory**.



Adapting CPU code to GPU code is not merely a matter of software engineering; it often requires the **redesign of the algorithm**

# Exascale Computing Project

- **Mission:** Tackle **real-world computational problems** with exascale computing.

Frontier @Oak Ridge



Aurora @Argonne



# Exascale Computing Project

- ▶ **Mission:** Tackle **real-world computational problems** with exascale computing.

Frontier @Oak Ridge



Aurora @Argonne



- ▶ **Goal:** Build a **comprehensive software infrastructure** for nonlinear optimization on GPUs.

# Exascale Computing Project

- **Mission:** Tackle **real-world computational problems** with exascale computing.

Frontier @Oak Ridge (AMD GPUs)



Aurora @Argonne (Intel GPUs)



- **Goal:** Build a **comprehensive software infrastructure** for nonlinear optimization on GPUs.
- **Challenge #1:** No software infrastructure for **classical nonlinear optimization** on GPUs.

# Exascale Computing Project

- ▶ **Mission:** Tackle **real-world computational problems** with exascale computing.

Frontier @Oak Ridge (AMD GPUs)



Aurora @Argonne (Intel GPUs)



- ▶ **Goal:** Build a **comprehensive software infrastructure** for nonlinear optimization on GPUs.
- ▶ **Challenge #1:** No software infrastructure for **classical nonlinear optimization** on GPUs.
- ▶ **Challenge #2:** **Heterogeneous** development environment (**NVIDIA**, **AMD**, and **Intel**).

# Exascale Computing Project

- ▶ **Mission:** Tackle **real-world computational problems** with exascale computing.

Frontier @Oak Ridge (AMD GPUs)



Aurora @Argonne (Intel GPUs)

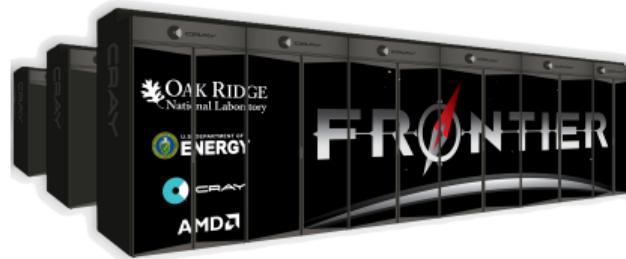


- ▶ **Goal:** Build a **comprehensive software infrastructure** for nonlinear optimization on GPUs.
- ▶ **Challenge #1:** No software infrastructure for **classical nonlinear optimization** on GPUs.
- ▶ **Challenge #2: Heterogeneous** development environment (**NVIDIA**, **AMD**, and **Intel**).
- ▶ Furthermore, we want to achieve
  - ▶ **Performance:** at least an order of magnitude speedup.

# Exascale Computing Project

- ▶ **Mission:** Tackle **real-world computational problems** with exascale computing.

Frontier @Oak Ridge (AMD GPUs)



Aurora @Argonne (Intel GPUs)

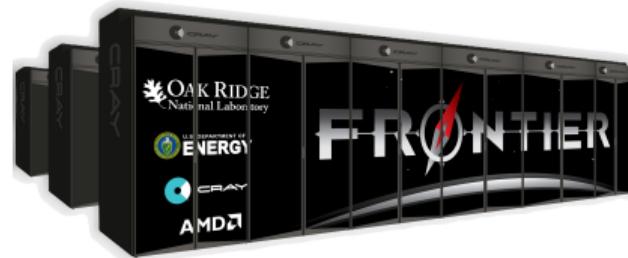


- ▶ **Goal:** Build a **comprehensive software infrastructure** for nonlinear optimization on GPUs.
- ▶ **Challenge #1:** No software infrastructure for **classical nonlinear optimization** on GPUs.
- ▶ **Challenge #2: Heterogeneous** development environment (**NVIDIA**, **AMD**, and **Intel**).
- ▶ Furthermore, we want to achieve
  - ▶ **Performance:** at least an order of magnitude speedup.
  - ▶ **Portability:** compatibility with **NVIDIA**, **AMD**, and **Intel**.

# Exascale Computing Project

- ▶ **Mission:** Tackle **real-world computational problems** with exascale computing.

Frontier @Oak Ridge (AMD GPUs)



Aurora @Argonne (Intel GPUs)



- ▶ **Goal:** Build a **comprehensive software infrastructure** for nonlinear optimization on GPUs.
- ▶ **Challenge #1:** No software infrastructure for **classical nonlinear optimization** on GPUs.
- ▶ **Challenge #2: Heterogeneous** development environment (**NVIDIA**, **AMD**, and **Intel**).
- ▶ Furthermore, we want to achieve
  - ▶ **Performance:** at least an order of magnitude speedup.
  - ▶ **Portability:** compatibility with **NVIDIA**, **AMD**, and **Intel**.
  - ▶ **Application:** energy infrastructure problems (AC optimal power flow, in particular).

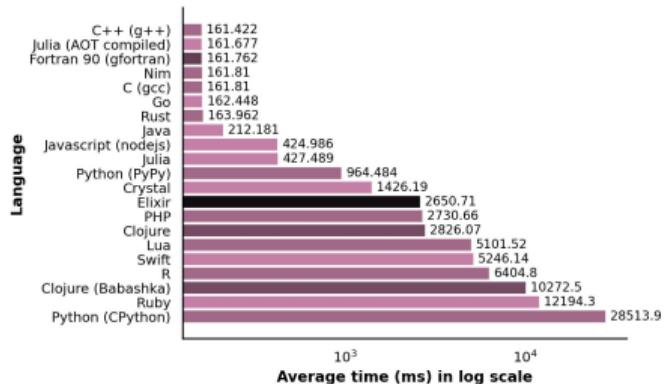
# Language of Choice: Julia



- Runs **as fast as C/C++/Fortran**, but is as simple as Python/R/MATLAB.

## Speed comparison of various programming languages

*Method: calculating  $\pi$  through the Leibniz formula 100000000 times*

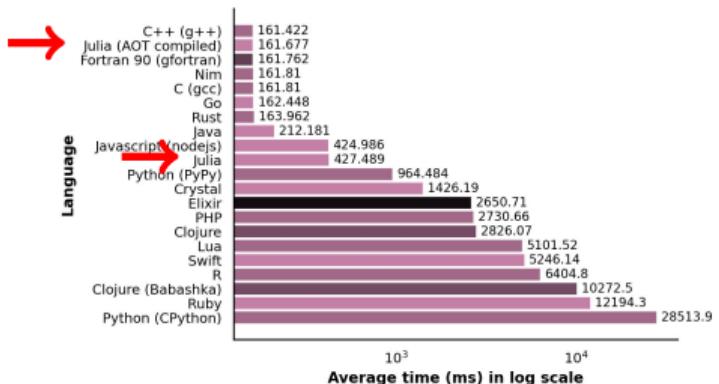


# Language of Choice: Julia

- ▶ Runs as fast as C/C++/Fortran, but is as simple as Python/R/MATLAB.

## Speed comparison of various programming languages

Method: calculating  $\pi$  through the Leibniz formula 100000000 times



Generated: 2022-10-16 19:55

<https://github.com/niklas-heer/speed-comparison>

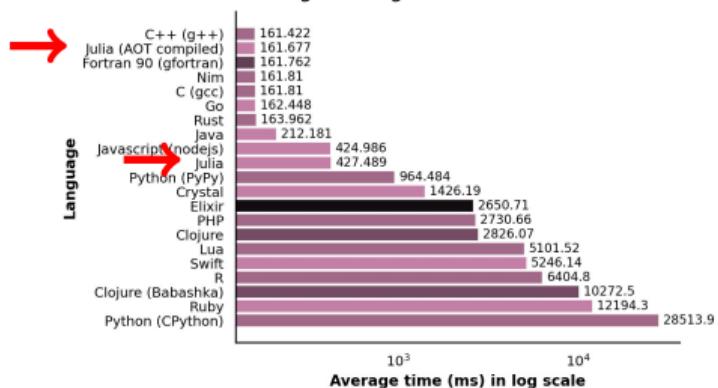
# Language of Choice: Julia



- ▶ Runs as fast as C/C++/Fortran, but is as simple as Python/R/MATLAB.

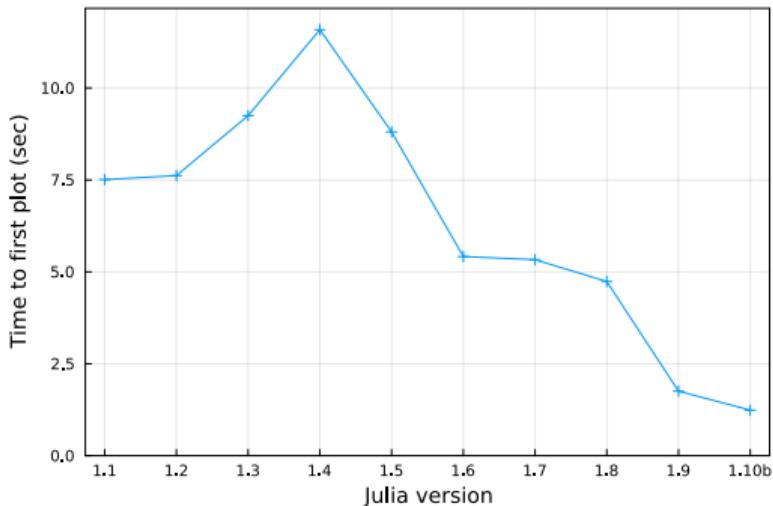
## Speed comparison of various programming languages

Method: calculating  $\pi$  through the Leibniz formula 100000000 times



Generated: 2022-10-16 19:55

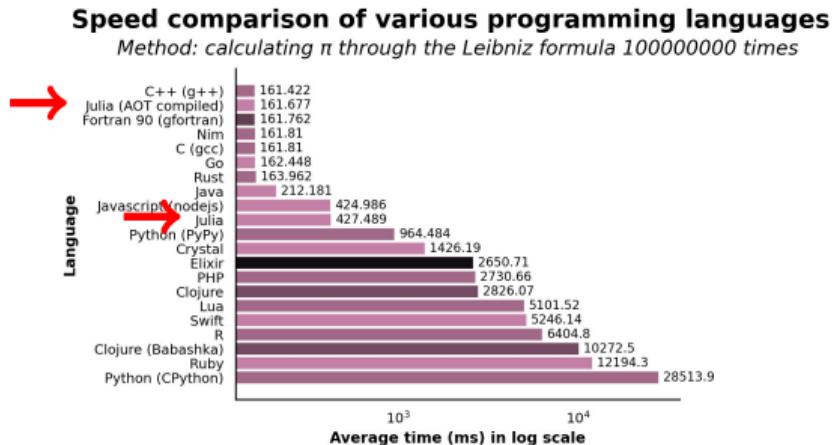
<https://github.com/niklas-heer/speed-comparison>



# Language of Choice: Julia

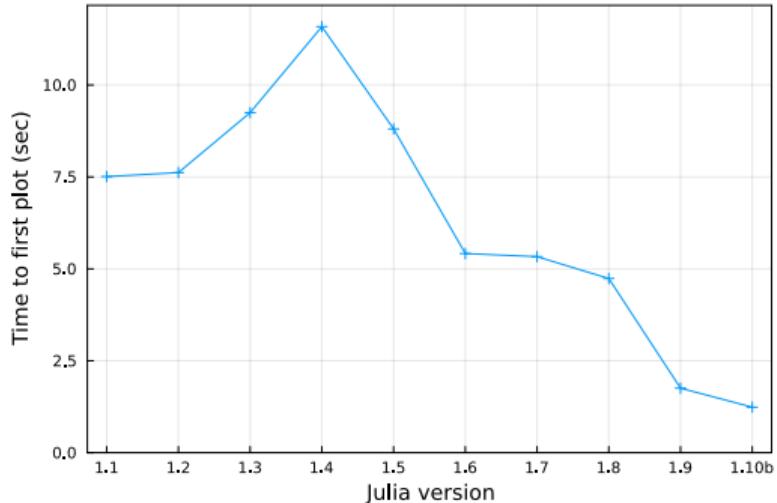


- ▶ Runs as fast as C/C++/Fortran, but is as simple as Python/R/MATLAB.



Generated: 2022-10-16 19:55

<https://github.com/niklas-heer/speed-comparison>



- ▶ Fast development: Julia resolves the “two-language problem”.

# Language of Choice: Julia



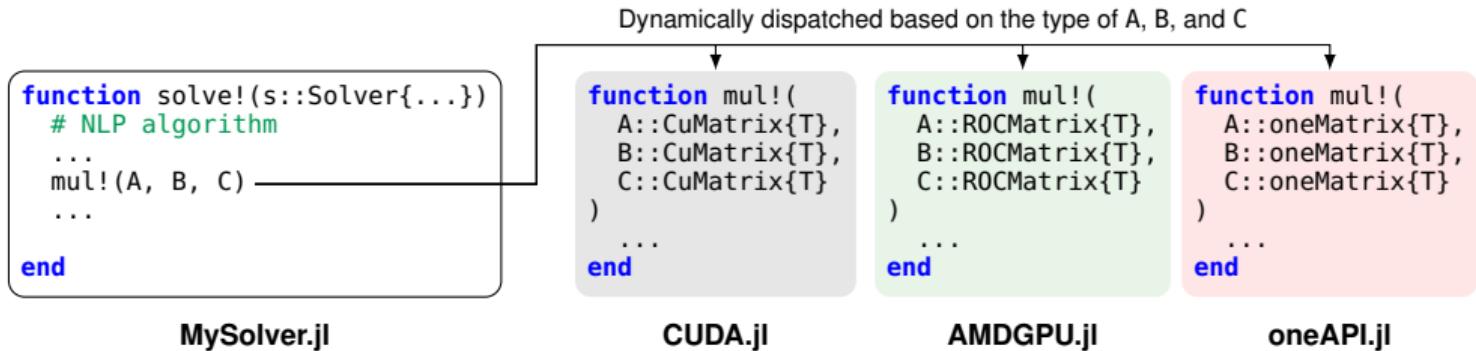
- ▶ **Multiple dispatch:** High-level abstraction while specializing for specific data types.

```
function solve!(s::Solver{...})
    # NLP algorithm
    ...
    mul!(A, B, C)
    ...
end
```

# Language of Choice: Julia



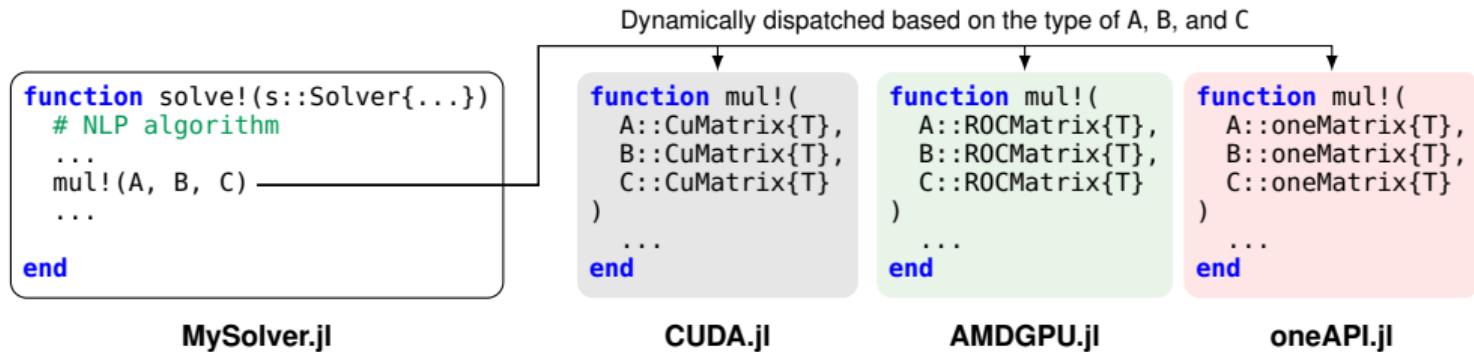
- ▶ **Multiple dispatch:** High-level abstraction while specializing for specific data types.



# Language of Choice: Julia



- ▶ **Multiple dispatch:** High-level abstraction while specializing for specific data types.



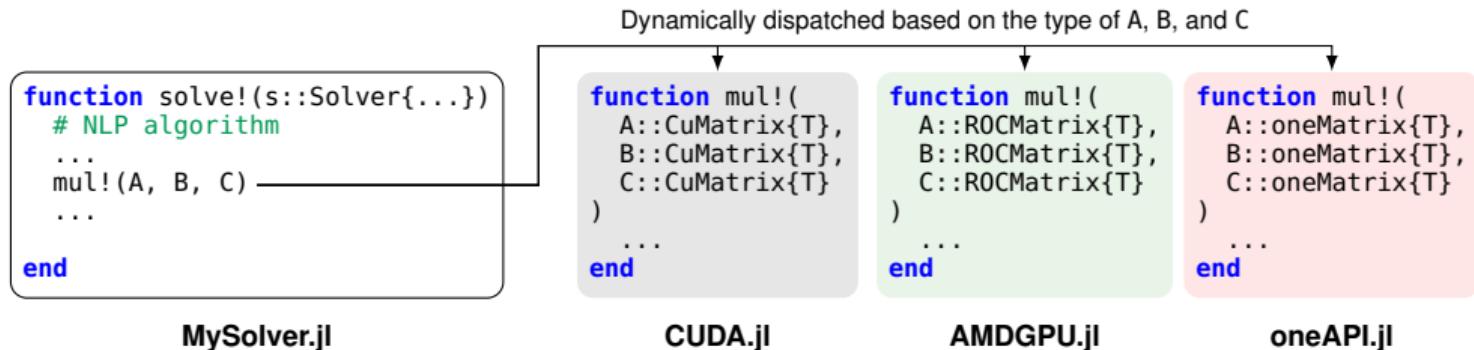
- ▶ **Portable kernel programming:** Compatibility across various architectures.

```
@kernel function _mul!(
    A,B,C
)
# portable GPU kernel
...
end
```

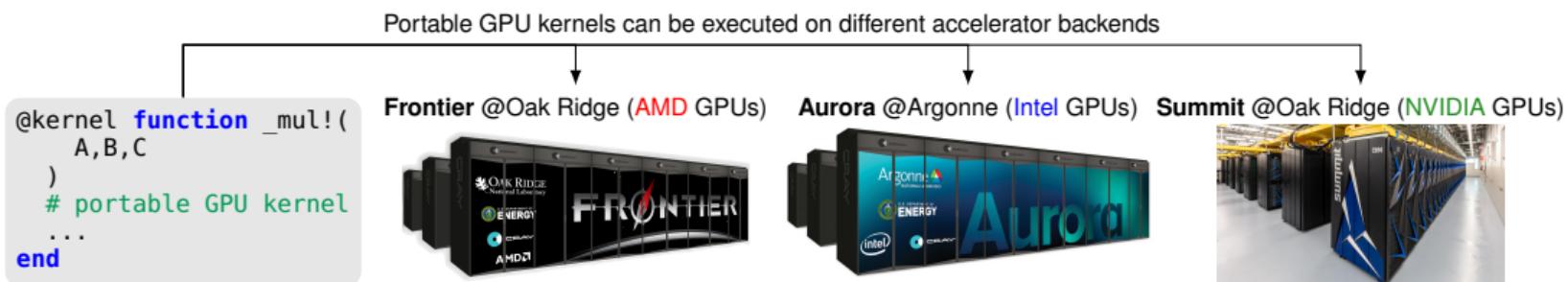
# Language of Choice: Julia



- ▶ **Multiple dispatch:** High-level abstraction while specializing for specific data types.



- ▶ **Portable kernel programming:** Compatibility across various architectures.



# Summary

GPU/HPC



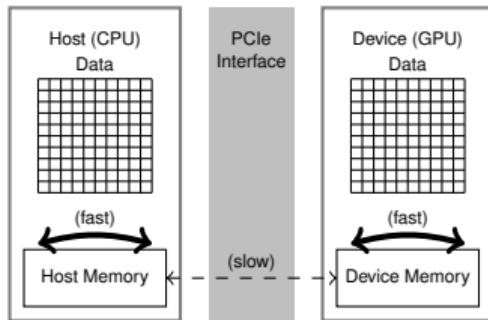
- ▶ Motivation: Harness accelerated computing for nonlinear programming.

# Summary



GPU/HPC

SIMD Architecture

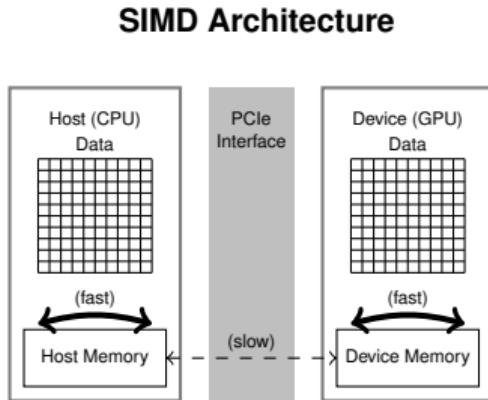


- ▶ **Motivation:** Harness accelerated computing for nonlinear programming.
- ▶ Adapting CPU code into GPU code is **not merely an issue of software engineering**.

# Summary



GPU/HPC



SIMD Architecture

Performance, Fast Development,  
& Portability



- ▶ **Motivation:** Harness accelerated computing for nonlinear programming.
- ▶ Adapting CPU code into GPU code is **not merely an issue of software engineering**.
- ▶ **Goal:** Build a **comprehensive software infrastructure** for nonlinear programming on GPUs on **Julia Language**, for its **fast performance, fast development speed**, and **portability**.

# Outline

1. Motivation
2. Nonlinear optimization framework
3. Implementing algebraic modeling systems for GPUs
4. Implementing nonlinear programming solvers for GPUs
  - 4.1. LiftedKKT
  - 4.2. HybridKKT
  - 4.3. Numerical results
5. Remaining challenges
6. Conclusion

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\min_{x \geq 0} f(x)$$

$$\text{s.t. } c(x) = 0$$

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\min_{x \geq 0} f(x)$$

$$\text{s.t. } c(x) = 0$$

- ▶ In classical problems (e.g., optimal power flow),
  - ▶ the objective and constraints are **smooth**
  - ▶ **large number of variables and constraints**
  - ▶ the problem is **highly sparse**.

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\min_{x \geq 0} f(x)$$

$$\text{s.t. } c(x) = 0$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

- ▶ In classical problems (e.g., optimal power flow),
  - ▶ the objective and constraints are **smooth**
  - ▶ **large number of variables and constraints**
  - ▶ the problem is **highly sparse**.
- ▶ Interior-point methods
  - ▶ Inequalities  $x \geq 0$  replaced by smooth log-barrier functions  $f(x) - \mu \sum_i \log(x[i])$ .
  - ▶ **Newton's Step** is computed by solving a “**KKT system**” (large, sparse, symmetric indefinite, ill-conditioned system).

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\min_{x \geq 0} f(x)$$

$$\text{s.t. } c(x) = 0$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$x^{(k+1)} = x^{(k)} + \alpha \Delta x$$

$$\lambda^{(k+1)} = \lambda^{(k)} + \alpha \Delta \lambda$$

- ▶ In classical problems (e.g., optimal power flow),
  - ▶ the objective and constraints are **smooth**
  - ▶ **large number of variables and constraints**
  - ▶ the problem is **highly sparse**.
- ▶ Interior-point methods
  - ▶ Inequalities  $x \geq 0$  replaced by smooth log-barrier functions  $f(x) - \mu \sum_i \log(x[i])$ .
  - ▶ **Newton's Step** is computed by solving a “**KKT system**” (large, sparse, symmetric indefinite, ill-conditioned system).
  - ▶ Line search (along with several additional heuristics) ensures **global convergence**.

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\begin{aligned} \min_{x \geq 0} f(x) \\ \text{s.t. } c(x) = 0 \end{aligned}$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha \Delta x \\ \lambda^{(k+1)} &= \lambda^{(k)} + \alpha \Delta \lambda \end{aligned}$$

## Algebraic Modeling Systems

AMPL, CasADi, JuMP, Gravity, ...



- ▶ Algebraic modeling systems provides **front-end** to specify models and (often) provides **derivative computation capabilities**.

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\begin{aligned} \min_{x \geq 0} f(x) \\ \text{s.t. } c(x) = 0 \end{aligned}$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha \Delta x \\ \lambda^{(k+1)} &= \lambda^{(k)} + \alpha \Delta \lambda \end{aligned}$$

## Algebraic Modeling Systems

AMPL, CasADi, JuMP, Gravity, ...

## Nonlinear Optimization Solvers

Ipopt, Knitro, Pynumero, ...

- ▶ Algebraic modeling systems provides **front-end** to specify models and (often) provides **derivative computation capabilities**.
- ▶ Nonlinear optimization solvers apply iterations of optimization algorithms.

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\begin{aligned} \min_{x \geq 0} f(x) \\ \text{s.t. } c(x) = 0 \end{aligned}$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha \Delta x \\ \lambda^{(k+1)} &= \lambda^{(k)} + \alpha \Delta \lambda \end{aligned}$$

## Algebraic Modeling Systems

AMPL, CasADi, JuMP, Gravity, ...

## Nonlinear Optimization Solvers

Ipopt, Knitro, Pynumero, ...

## Sparse Linear Solvers

HSL (ma27, ma57, ...), Pardiso, ...

- ▶ Algebraic modeling systems provides **front-end** to specify models and (often) provides **derivative computation capabilities**.
- ▶ Nonlinear optimization solvers apply iterations of optimization algorithms.
- ▶ Sparse linear solvers solves KKT systems using **sparse matrix factorization**.

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\begin{aligned} \min_{x \geq 0} f(x) \\ \text{s.t. } c(x) = 0 \end{aligned}$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha \Delta x \\ \lambda^{(k+1)} &= \lambda^{(k)} + \alpha \Delta \lambda \end{aligned}$$

## Algebraic Modeling Systems

AMPL, CasADi, JuMP, Gravity, ...

## Nonlinear Optimization Solvers

Ipopt, Knitro, Pynumero, ...

## Sparse Linear Solvers

HSL (ma27, ma57, ...), Pardiso, ...

- ▶ These software tools have enabled the success of nonlinear optimization on CPUs.

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\begin{aligned} \min_{x \geq 0} f(x) \\ \text{s.t. } c(x) = 0 \end{aligned}$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha \Delta x \\ \lambda^{(k+1)} &= \lambda^{(k)} + \alpha \Delta \lambda \end{aligned}$$

## Algebraic Modeling Systems

AMPL, CasADi, JuMP, Gravity, ...

## Nonlinear Optimization Solvers

Ipopt, Knitro, Pynumero, ...

## Sparse Linear Solvers

HSL (ma27, ma57, ...), Pardiso, ...

- ▶ These software tools have enabled the success of nonlinear optimization on CPUs.
- ▶ Many software tools have been developed in 80s-90s (**heavily optimized for CPUs**).

# Nonlinear optimization framework: state-of-the-art on CPU

## Problem Formulation

$$\begin{aligned} \min_{x \geq 0} f(x) \\ \text{s.t. } c(x) = 0 \end{aligned}$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha \Delta x \\ \lambda^{(k+1)} &= \lambda^{(k)} + \alpha \Delta \lambda \end{aligned}$$

## Algebraic Modeling Systems

AMPL, CasADi, JuMP, Gravity, ...

## Nonlinear Optimization Solvers

Ipopt, Knitro, Pynumero, ...

## Sparse Linear Solvers

HSL (ma27, ma57, ...), Pardiso, ...



- ▶ These software tools have enabled the success of nonlinear optimization on CPUs.
- ▶ Many software tools have been developed in 80s-90s (**heavily optimized for CPUs**).
- ▶ Now we need **GPU-equivalent** of these tools.

# Nonlinear optimization framework on GPUs

## Problem Formulation

$$\begin{aligned} \min_{x \geq 0} f(x) \\ \text{s.t. } c(x) = 0 \end{aligned}$$

## Newton's Step Computation

$$\underbrace{\begin{bmatrix} W + \Sigma + \delta_w I & A^\top \\ A & -\delta_c I \end{bmatrix}}_{\text{"KKT System" (ill-conditioned)}} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} p^x \\ p^\lambda \end{bmatrix}$$

## Line Search

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha \Delta x \\ \lambda^{(k+1)} &= \lambda^{(k)} + \alpha \Delta \lambda \end{aligned}$$

Algebraic Modeling Systems

Nonlinear Optimization Solvers

Sparse Linear Solvers



ExaModels

- Parallel AD with SIMD abstraction
- Runs on GPU architectures

<https://github.com/exanauts/ExaModels.jl>  
<https://github.com/MadNLP/MadNLP.jl>  
<https://github.com/exanauts/CUDSS.jl>



MadNLP

- Lifted & Hybrid KKT System
- Runs on NVIDIA GPUs



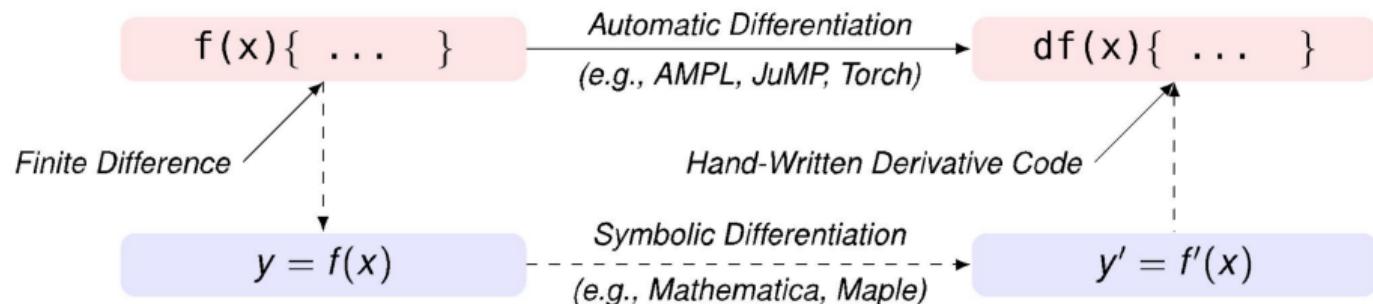
- Parallel Cholesky factorization
- Runs on NVIDIA GPUs

# Outline

1. Motivation
2. Nonlinear optimization framework
3. Implementing algebraic modeling systems for GPUs
4. Implementing nonlinear programming solvers for GPUs
  - 4.1. LiftedKKT
  - 4.2. HybridKKT
  - 4.3. Numerical results
5. Remaining challenges
6. Conclusion

# Sparse Automatic Differentiation (AD)

- Different paradigms for derivative computation:



- AD is superior in terms of **efficiency**, **accuracy**, and **convenience**
- AD has been the **state-of-the-art** method in numerical optimization
- **Parallelizing** sparse automatic differentiation on GPUs is **non-trivial** due to **irregularities**

# SIMD Abstraction



- Large-scale optimization problems almost always have repeated patterns
- SIMD Abstraction can capture such repeated patterns:

$$\begin{aligned} \min_{x^b \leq x \leq x^u} & \sum_{l \in [L]} \sum_{i \in [I_l]} f^{(l)}(x; p_i^{(l)}) \\ \text{s.t. } & \sum_{n \in [N]} \sum_{k \in [K_n]} h^{(n)}(x; s_k^{(n)}) = 0 \end{aligned}$$

Annotations:

- a small number of different patterns (points to the first term in the objective)
- "single instruction" (points to the inner sum in the objective)
- "over multiple data" (points to the outer sum in the constraint)

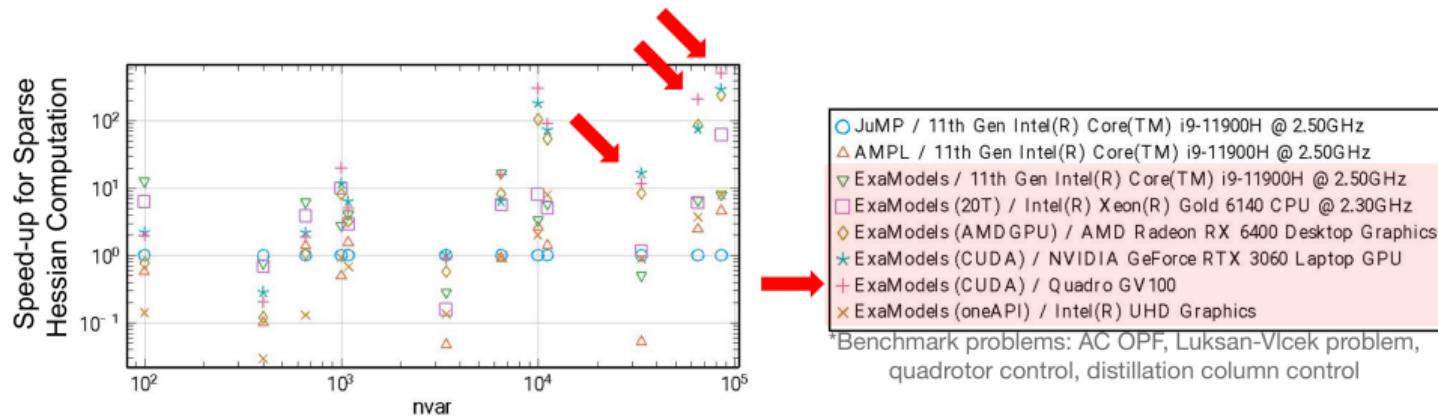
- Repeated patterns are inputted as **iterators** (data can be stored in structured format)

```
constraint(c, 3 * x[i+1]^3 + 2 * sin(x[i+2]) for i = 1:N-2)
```



- For each pattern, the AD kernel is **compiled** and **executed over multiple data** in parallel

# Sparse AD Benchmark



- For the largest case, **ExaModels on GPU** is **100x faster** than the state-of-the-art tools on CPUs
- ExaModels runs on **all major GPU architectures** and **single/multi-threaded CPUs**

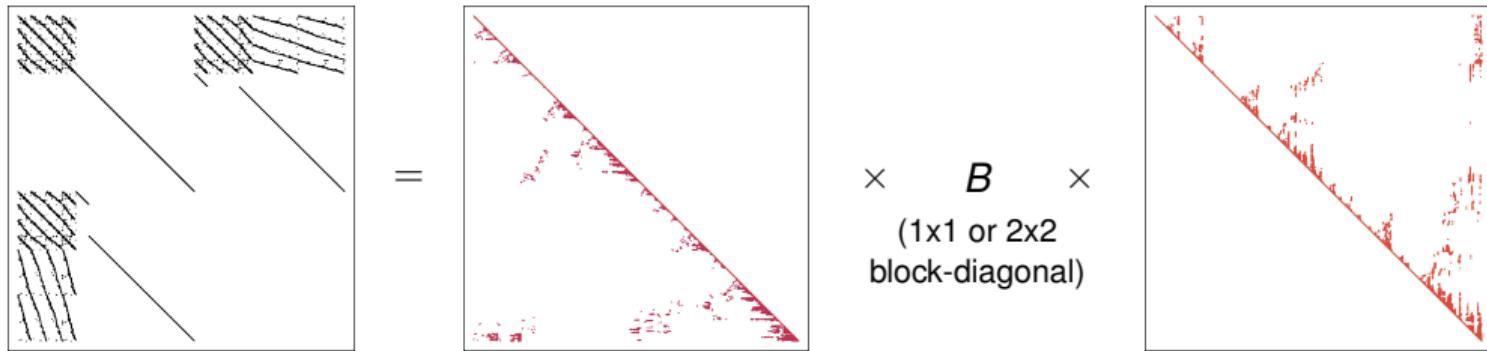
**Sparse AD with SIMD abstraction enables efficient derivative computations on GPUs**

# Outline

1. Motivation
2. Nonlinear optimization framework
3. Implementing algebraic modeling systems for GPUs
4. Implementing nonlinear programming solvers for GPUs
  - 4.1. LiftedKKT
  - 4.2. HybridKKT
  - 4.3. Numerical results
5. Remaining challenges
6. Conclusion

# Why it is Challenging: Sparse Linear Solvers

- Solving KKT systems on CPUs has traditionally relied on **direct  $LBL^\top$  factorization**.



- $LBL^\top$  factorization requires **numerical pivoting**, which is **challenging to parallelize** [Swirydowicz, 2024].
- First-order methods [Lu, 2023] and iterative solvers [Schubiger, 2020] [Cao, 2016] have been proposed but can be ineffective/less reliable for nonconvex problems.
- We aim to develop a "pivoting-free" interior-point method.

# LiftedKKT

- We aim to transform the KKT system into an equivalent **positive definite system**, which can be **factorized without pivoting** based on **Choleseky factorization** (solver available for NVIDIA)
- This can be achieved by:
  - Lift and relax** the equalities into **inequalities**:  $g(x) = 0 \implies g(x) - s = 0$ ,  $s^b \leq s \leq s^{\#}$ ,
  - Eliminate the slack variables** (so-called condensation):

$$\begin{array}{c} \left[ \begin{array}{cccccc} W^{(\ell)} + \delta_w^{(\ell)} I & A^{(\ell)\top} & -I & I & -I & I \\ \delta_w^{(\ell)} I & -I & -\delta_c^{(\ell)} I & X^{(\ell)} - X^b & & \\ A^{(\ell)} & -I & -\delta_c^{(\ell)} I & & & \\ Z_x^{(\ell)b} & X^{(\ell)} - X^b & & & & \\ -Z_x^{(\ell)\#} & & & S^{(\ell)} - S^b & & \\ Z_s^{(\ell)b} & & & & S^{\#} - S^{(\ell)} & \\ -Z_s^{(\ell)\#} & & & & & \end{array} \right] \begin{bmatrix} \Delta x \\ \Delta s \\ \Delta y \\ \Delta z_x^b \\ \Delta z_x^{\#} \\ \Delta z_s^b \\ \Delta z_s^{\#} \\ \Delta z_s^{\#} \end{bmatrix} = \begin{bmatrix} p_x^{(\ell)} \\ p_s^{(\ell)} \\ p_y^{(\ell)} \\ p_z_x^{(\ell)} \\ p_z_x^{(\ell)} \\ p_z_s^{(\ell)} \\ p_z_s^{(\ell)} \\ p_z_s^{(\ell)} \end{bmatrix} \\ \text{Nonsingular} \quad \implies (W + \delta_w I + \Sigma_x + A^\top D A) \Delta x = q_x + A^\top (C q_s + D q_y) \quad \text{Positive-definite} \end{array}$$

- Condensation** causes **extreme but relatively benign ill-conditioning** [M. Wright, 1998]

# HybridKKT

- Direct + iterative approach, with **regularized equalities**

- Condense KKT systems by **eliminating inequalities**:

$$\begin{bmatrix} W_k + \delta_x I & H_k^\top & G_k^\top \\ D_s + \delta_x I & I & \\ H_k & I & -\delta_c I \\ G_k & & -\delta_c I \end{bmatrix} \begin{bmatrix} d_x \\ d_s \\ d_z \\ d_y \end{bmatrix} = - \begin{bmatrix} r_1 \\ r_2 \\ r_4 \\ r_3 \end{bmatrix} \longrightarrow \underbrace{\begin{bmatrix} K_k & G_k^\top \\ G_k & 0 \end{bmatrix}}_{\substack{K_{cond} \\ \text{Nonsingular}}} \begin{bmatrix} d_x \\ d_y \end{bmatrix} = - \begin{bmatrix} r_1 + H_k^\top(D_s r_4 - r_2) \\ r_3 \end{bmatrix}$$

Still indefinite

- Regularize equalities with sufficiently large  $\gamma$  :

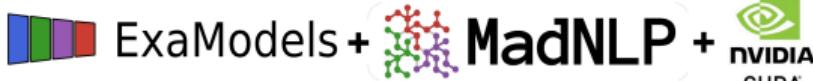
$$\text{Positive-definite} \leftarrow K_\gamma := \begin{bmatrix} K_k + \gamma G_k^\top G_k & G_k^\top \\ G_k & 0 \end{bmatrix} \begin{bmatrix} d_x \\ d_y \end{bmatrix} = - \begin{bmatrix} r_\gamma \\ r_3 \end{bmatrix}$$

- Solve the Schur complement system with iterative method (e.g., conjugate gradient):

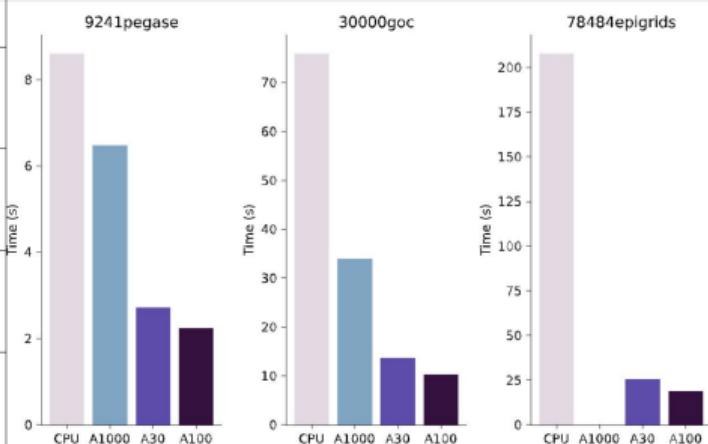
Converges to  $I$  as  $\gamma \rightarrow \infty$   $\leftarrow (G_k K_\gamma^{-1} G_k^\top) d_y = r_3 - G_k K_\gamma^{-1} r_\gamma$

- Similarly to the “Lifted KKT System”, we have **extreme but relatively benign ill-conditioning**

# AC Optimal Power Flow



Case	CPU				Lifted KKT on GPU				Hybrid KKT on GPU			
	it	init	lin	total	it	init	lin	total	it	init	lin	total
89_pegase	32	0.00	0.02	<b>0.03</b>	29	0.03	0.12	<b>0.24</b>	32	0.03	0.07	<b>0.22</b>
179_goc	45	0.00	0.03	<b>0.05</b>	39	0.03	0.19	<b>0.35</b>	45	0.03	0.07	<b>0.25</b>
500_goc	39	0.01	0.10	<b>0.14</b>	39	0.05	0.09	<b>0.26</b>	39	0.05	0.07	<b>0.27</b>
793_goc	35	0.01	0.12	<b>0.18</b>	57	0.06	0.27	<b>0.52</b>	35	0.05	0.10	<b>0.30</b>
1354_pegase	49	0.02	0.35	<b>0.52</b>	96	0.12	0.69	<b>1.22</b>	49	0.12	0.17	<b>0.50</b>
2000_goc	42	0.03	<b>0.66</b>	<b>0.93</b>	46	0.15	0.30	<b>0.66</b>	42	0.16	0.14	<b>0.50</b>
2312_goc	43	0.02	0.59	<b>0.82</b>	45	0.11	0.32	<b>0.68</b>	43	0.14	0.21	<b>0.56</b>
2742_goc	125	0.04	3.76	<b>7.31</b>	157	0.20	1.93	<b>15.49</b>	-	-	-	-
2869_pegase	55	0.04	1.09	<b>1.52</b>	57	0.20	0.30	<b>0.80</b>	55	0.21	0.26	<b>0.73</b>
3022_goc	55	0.03	0.98	<b>1.39</b>	48	0.18	0.23	<b>0.66</b>	55	0.18	0.23	<b>0.68</b>
3970_goc	48	0.05	1.95	<b>2.53</b>	47	0.26	0.37	<b>0.87</b>	48	0.27	0.24	<b>0.80</b>
4020_goc	59	0.06	3.90	<b>4.60</b>	123	0.28	1.75	<b>3.15</b>	59	0.29	0.41	<b>1.08</b>
4601_goc	71	0.09	3.09	<b>4.16</b>	67	0.27	0.51	<b>1.17</b>	71	0.28	0.39	<b>1.12</b>
4619_goc	49	0.07	3.21	<b>3.91</b>	49	0.31	0.59	<b>1.25</b>	49	0.33	0.31	<b>0.95</b>
4837_goc	59	0.08	2.49	<b>3.33</b>	59	0.29	0.58	<b>1.31</b>	59	0.29	0.35	<b>0.98</b>
4917_goc	63	0.07	1.97	<b>2.72</b>	55	0.26	0.55	<b>1.18</b>	63	0.36	0.34	<b>0.94</b>
5658_epigrids	51	0.31	2.80	<b>3.86</b>	58	0.35	0.66	<b>1.51</b>	51	0.35	0.35	<b>1.03</b>
7336_epigrids	50	0.13	3.60	<b>4.91</b>	56	0.45	0.95	<b>1.89</b>	50	0.43	0.35	<b>1.13</b>
8387_pegase	74	0.14	5.31	<b>7.62</b>	82	0.59	0.79	<b>2.30</b>	75	0.58	0.66	<b>8.84</b>
9241_pegase	74	0.15	6.11	<b>8.60</b>	101	0.63	0.88	<b>2.76</b>	71	0.63	0.99	<b>2.24</b>
9591_goc	67	0.20	11.14	<b>13.37</b>	98	0.63	2.67	<b>4.58</b>	67	0.62	0.74	<b>1.96</b>
10000_goc	82	0.15	6.00	<b>8.16</b>	64	0.49	0.81	<b>1.83</b>	82	0.49	0.75	<b>1.82</b>
10192_epigrids	54	0.41	7.79	<b>10.08</b>	57	0.67	1.14	<b>2.40</b>	54	0.67	0.66	<b>1.81</b>
10180_goc	71	0.24	12.04	<b>14.74</b>	67	0.75	0.99	<b>2.72</b>	71	0.71	1.09	<b>2.50</b>
13659_pegase	63	0.45	7.21	<b>10.14</b>	75	0.83	1.05	<b>2.96</b>	62	0.81	0.93	<b>2.47</b>
19402_goc	69	0.63	31.71	<b>36.92</b>	73	1.42	2.28	<b>5.38</b>	69	1.44	1.93	<b>4.31</b>
20758_epigrids	51	0.63	14.27	<b>18.21</b>	53	1.34	1.05	<b>3.57</b>	51	1.35	1.55	<b>3.51</b>
30000_goc	183	0.65	63.02	<b>75.95</b>	-	-	-	-	225	1.22	5.59	<b>10.27</b>
78484_epigrids	102	2.57	179.29	<b>207.79</b>	101	5.94	5.62	<b>18.03</b>	104	6.29	9.01	<b>18.90</b>



Optimizing entire eastern interconnection

Table 3 OPF benchmark, solved with a tolerance  $\text{tol}=1\text{e-}6$ . (A100 GPU)

- For large-scale cases ( $> 20k$  vars), GPU becomes significantly faster than CPU (up to  $\times 10$ )

# Distillation Column



#time steps	CPU			Lifted KKT on GPU			Hybrid KKT on GPU		
	init (s)	it	solve (s)	init (s)	it	solve (s)	init (s)	it	solve (s)
100	0.1	7	0.1	0.1	11	0.1	0.1	7	0.0
500	0.1	7	0.5	0.2	12	0.1	0.2	7	0.1
1,000	0.1	7	1.5	0.4	12	0.2	0.4	7	0.1
5,000	0.6	7	8.2	2.3	13	0.5	2.3	7	0.4
10,000	1.3	7	18.7	5.2	13	0.9	5.3	7	0.7
20,000	4.3	7	38.2	10.7	14	2.1	11.3	7	1.4
50,000	15.9	7	98.8	30.4	14	5.5	31.2	7	3.8

- “Symbolic analysis” is often the bottleneck on GPUs, but this can be computed “off-line”  
thus, **online computation performance** can be **even greater**
- The distillation column control problem can be solved more than 20x faster

**ExaModels, MadNLP, and CUDSS provide efficient and reliable solution framework for large-scale nonlinear optimization problems**

# Outline

1. Motivation
2. Nonlinear optimization framework
3. Implementing algebraic modeling systems for GPUs
4. Implementing nonlinear programming solvers for GPUs
  - 4.1. LiftedKKT
  - 4.2. HybridKKT
  - 4.3. Numerical results
5. Remaining challenges
6. Conclusion

## Accuracy

- We have **avoided indefinite systems** by replacing them with **positive definite systems**.

# Accuracy

- ▶ We have **avoided indefinite systems** by replacing them with **positive definite systems**.
- ▶ Challenges arise from **ill-conditioning of condensed KKT systems**.

$$\begin{bmatrix} H + \Sigma_x & J^\top \\ \Sigma_s & -I \\ J & -I \end{bmatrix} \xrightarrow{\text{condensation}} H + \Sigma_x + J^\top \Sigma_s^{-1} J$$

# Accuracy

- ▶ We have **avoided indefinite systems** by replacing them with **positive definite systems**.
- ▶ Challenges arise from **ill-conditioning of condensed KKT systems**.

$$\begin{bmatrix} H + \Sigma_x & J^\top \\ \Sigma_s & -I \\ J & -I \end{bmatrix} \xrightarrow{\text{condensation}} H + \Sigma_x + J^\top \Sigma_s^{-1} J$$

Causes ill-conditioning

# Accuracy

- ▶ We have **avoided indefinite systems** by replacing them with **positive definite systems**.
- ▶ Challenges arise from **ill-conditioning of condensed KKT systems**.

$$\begin{bmatrix} H + \Sigma_x & J^\top \\ \Sigma_s & -I \\ J & -I \end{bmatrix} \xrightarrow{\text{condensation}} H + \Sigma_x + J^\top \Sigma_s^{-1} J$$

Causes ill-conditioning

- ▶ **Alternative methods**, such as penalty method or augmented Lagrangian, **rely on similar manipulations**.

# Accuracy

- ▶ We have **avoided indefinite systems** by replacing them with **positive definite systems**.
- ▶ Challenges arise from **ill-conditioning of condensed KKT systems**.

$$\begin{bmatrix} H + \Sigma_x & J^\top \\ \Sigma_s & -I \\ J & -I \end{bmatrix} \xrightarrow{\text{condensation}} H + \Sigma_x + J^\top \Sigma_s^{-1} J$$

Causes ill-conditioning

- ▶ **Alternative methods**, such as penalty method or augmented Lagrangian, **rely on similar manipulations**.
- ▶ Work-arounds?
  - ▶ Using CPUs as we approach the solution.
  - ▶ Quadruple precision—challenging due to low-level kernel supports.
  - ▶ Hopefully, some breakthrough in sparse linear algebra (e.g., scalable parallel pivoting).

# Portability

Table: GPU Compatibility of Nonlinear Optimization Frameworks

		CPU (single)	CPU (multi)	NVIDIA GPU	AMD GPU	Intel GPU	Apple Metal
Algebraic Modeling Platforms	<b>AMPL</b>	✓	X	X	X	X	X
	<b>JuMP</b>	✓	X	X	X	X	X
	<b>ExaModels</b>	✓	✓	✓	✓	✓	X
NLP Solvers	<b>Ipopt</b>	✓	X	X	X	X	X
	<b>MadNLP</b>	✓	X	✓	X	X	X

- ▶ **ExaModels** has **full compatibility** with multi-threaded CPUs, **NVIDIA**, **AMD**, and **Intel** GPUs.

# Portability

Table: GPU Compatibility of Nonlinear Optimization Frameworks

		CPU (single)	CPU (multi)	NVIDIA GPU	AMD GPU	Intel GPU	Apple Metal
Algebraic Modeling Platforms	<b>AMPL</b>	✓	X	X	X	X	X
	<b>JuMP</b>	✓	X	X	X	X	X
	<b>ExaModels</b>	✓	✓	✓	✓	✓	X
NLP Solvers	<b>Ipopt</b>	✓	X	X	X	X	X
	<b>MadNLP</b>	✓	X	✓	X	X	X

- ▶ **ExaModels** has **full compatibility** with multi-threaded CPUs, **NVIDIA**, **AMD**, and **Intel** GPUs.
- ▶ **MadNLP** is currently **only compatible with NVIDIA**,

# Portability

Table: GPU Compatibility of Nonlinear Optimization Frameworks

		CPU (single)	CPU (multi)	NVIDIA GPU	AMD GPU	Intel GPU	Apple Metal
Algebraic Modeling Platforms	<b>AMPL</b>	✓	X	X	X	X	X
	<b>JuMP</b>	✓	X	X	X	X	X
	<b>ExaModels</b>	✓	✓	✓	✓	✓	X
NLP Solvers	<b>Ipopt</b>	✓	X	X	X	X	X
	<b>MadNLP</b>	✓	X	✓	X	X	X

- ▶ **ExaModels** has **full compatibility** with multi-threaded CPUs, **NVIDIA**, **AMD**, and **Intel** GPUs.
- ▶ **MadNLP** is currently **only compatible with NVIDIA**, but making it portable is not difficult.

# Portability

Table: GPU Compatibility of Nonlinear Optimization Frameworks

		CPU (single)	CPU (multi)	NVIDIA GPU	AMD GPU	Intel GPU	Apple Metal
Algebraic Modeling Platforms	<b>AMPL</b>	✓	X	X	X	X	X
	<b>JUMP</b>	✓	X	X	X	X	X
	<b>ExaModels</b>	✓	✓	✓	✓	✓	X
NLP Solvers	<b>Ipopt</b>	✓	X	X	X	X	X
	<b>MadNLP</b>	✓	X	✓	X	X	X

- ▶ **ExaModels** has **full compatibility** with multi-threaded CPUs, **NVIDIA**, **AMD**, and **Intel** GPUs.
- ▶ **MadNLP** is currently **only compatible with NVIDIA**, but making it portable is not difficult.
- ▶ **Main obstacle:** AMD and Intel GPUs are limited in terms of sparse linear solvers.

# Portability

Table: GPU Compatibility of Nonlinear Optimization Frameworks

		CPU (single)	CPU (multi)	NVIDIA GPU	AMD GPU	Intel GPU	Apple Metal
Algebraic Modeling Platforms	<b>AMPL</b>	✓	X	X	X	X	X
	<b>JUMP</b>	✓	X	X	X	X	X
	<b>ExaModels</b>	✓	✓	✓	✓	✓	X
NLP Solvers	<b>Ipopt</b>	✓	X	X	X	X	X
	<b>MadNLP</b>	✓	X	✓	X	X	X

- ▶ **ExaModels** has **full compatibility** with multi-threaded CPUs, **NVIDIA**, **AMD**, and **Intel** GPUs.
- ▶ **MadNLP** is currently **only compatible with NVIDIA**, but making it portable is not difficult.
- ▶ **Main obstacle:** AMD and Intel GPUs are limited in terms of sparse linear solvers.
- ▶ Work-arounds?
  - ▶ Preconditioned iterative solvers.
  - ▶ Development of sparse linear solvers on GPU architectures?  
Work-in-Progress: A.Montoison and I. Duff. *MA57 on GPU !!!*

# Outline

1. Motivation
2. Nonlinear optimization framework
3. Implementing algebraic modeling systems for GPUs
4. Implementing nonlinear programming solvers for GPUs
  - 4.1. LiftedKKT
  - 4.2. HybridKKT
  - 4.3. Numerical results
5. Remaining challenges
6. Conclusion

# Conclusion

- ▶ Optimization on GPUs is a growing area (LPs, QPs, domain-specific problems)!

# Conclusion

- ▶ Optimization on GPUs is a growing area (LPs, QPs, domain-specific problems)!
- ▶ **GPU hardware** offers significant potential for **accelerating large-scale optimization**.

## Conclusion

- ▶ Optimization on GPUs is a growing area (LPs, QPs, domain-specific problems)!
- ▶ **GPU hardware** offers significant potential for **accelerating large-scale optimization**.
- ▶ Porting algorithms on GPUs often **requires complete redesign of the algorithms**.

## Conclusion

- ▶ Optimization on GPUs is a growing area (LPs, QPs, domain-specific problems)!
- ▶ **GPU hardware** offers significant potential for **accelerating large-scale optimization**.
- ▶ Porting algorithms on GPUs often **requires complete redesign of the algorithms**.
- ▶ We have achieved promising results: up to **20x faster solutions with moderate accuracy**.

## Conclusion

- ▶ Optimization on GPUs is a growing area (LPs, QPs, domain-specific problems)!
- ▶ **GPU hardware** offers significant potential for **accelerating large-scale optimization**.
- ▶ Porting algorithms on GPUs often **requires complete redesign of the algorithms**.
- ▶ We have achieved promising results: up to **20x faster solutions with moderate accuracy**.
- ▶ Challenges remain: **Ill-conditioning of condensed KKT system** and **portability**.

# Conclusion

- ▶ Optimization on GPUs is a growing area (LPs, QPs, domain-specific problems)!
- ▶ **GPU hardware** offers significant potential for **accelerating large-scale optimization**.
- ▶ Porting algorithms on GPUs often **requires complete redesign of the algorithms**.
- ▶ We have achieved promising results: up to **20x faster solutions with moderate accuracy**.
- ▶ Challenges remain: **Ill-conditioning of condensed KKT system** and **portability**.
- ▶ We envision **expanding the application scope of nonlinear programming**.
  - ▶ **Extremely large-scale** problems (coupled infrastructures, multi-stage, multiscale).
  - ▶ Problems involving **expensive surrogate models** (neural nets, simulations).