# IMPERIAL

# Learning Bounds on Computational Values in Iterative Methods using Reachability Analysis

**Mukund Verma,**
**Dr Ian McInerney,**
**Dr Ludovic Renson**

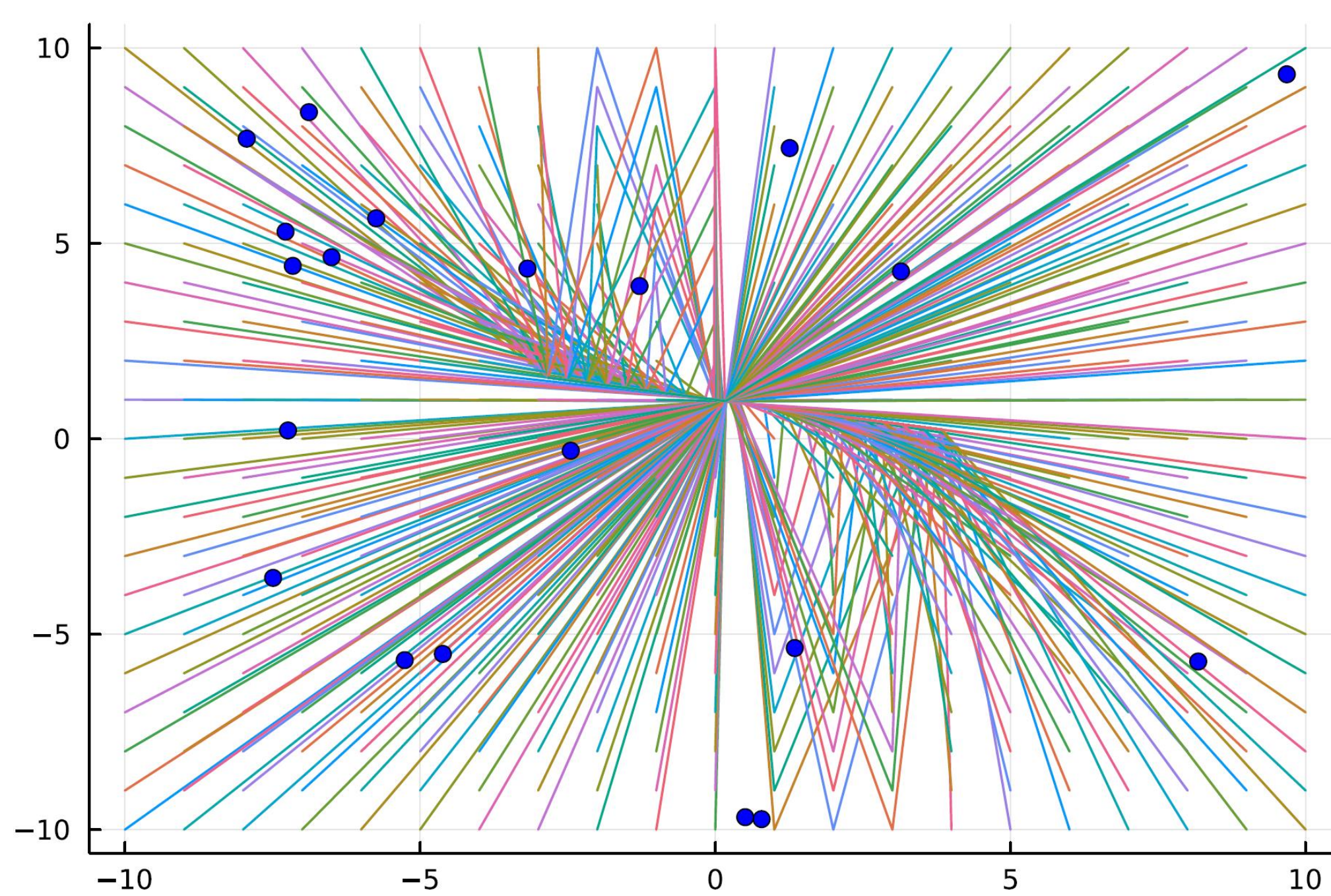**Dept. of Mechanical Engineering**

## Rationale

Each new generation of GPUs contain smaller number formats with a smaller *representable range*. Analysing how algorithms will behave with these smaller formats is essential, since algorithms that previously worked in higher precision may fail to converge to the correct result or become inaccurate when using these smaller number formats. However, the majority of the prior numerical analysis of algorithms in floating-point has ignored the representable range, since the formats found in modern computers have a representable range which is large enough to not affect most computations.

## Proposed Method

We propose analyzing the representable range necessary to implement specific algorithms by treating the algorithm as a non-linear discrete-time dynamical system and then performing reachability analysis on a Koopman operator linearization of the algorithm. Using the data-driven Koopman operator estimation for this analysis allows for building the Koopman operator of complex, or even black-box, numerical codes by only capturing the necessary algorithm state information at each iteration. The learned Koopman operator is then used to compute the reachable sets of the iterative method, which then provide the relevant information to determine what data types will have a suitable representable range.

## Illustrative Example

### Setup

- Gauss-Seidel stationary iterative method to solve $Ax = b$ with $A = \begin{bmatrix} 6 & 2 \\ 1 & 5 \end{bmatrix}, b = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$
- Use 20 Gaussian Radial Basis Functions with randomly placed centers and original system states as observables
- Construct Koopman operator using trajectory samples in the space [-10, 10] x [-10, 10]
- Approximate reachable set using Monte Carlo simulation of learned Koopman operator's trajectories

### Results

- Identify $x_2 \geq 0$ as an invariant set – stays positive if the number starts positive
- Approximated reachable invariant sets can inform required number format - $x_2$ could be unsigned number while $x_1$ must be a signed number

### Sampled trajectories and RBF centers



## Further work

This work is very preliminary and there are many open questions and possible research directions to improve the method and analysis:

- How to choose the observables for complicated non-linear/blackbox algorithms – learned observables?
- Can the algorithm's Koopman operator generalize to classes of problems?
- Formal guarantees in the reachability analysis
  - Use zonotopes in observable space
  - Reachability directly during Koopman operator estimation
- Integrating this method into a design or co-design workflow to optimize the number format or numerical method

### Simulated reachability from initial conditions using the Koopman operator



[-10, 10] x [0, 10]

[-10, 10] x [-10, 5]

[5, 10] x [-10, -5]

[0, 10] x [0, 10]

Legend: Initial set, Iteration 1, Iteration 2, Iteration 3, Iteration 4, Iteration 5, Iteration 6, Solution