

Transprecision Techniques for Linear Solvers and Non-Linear Regressions

JunKyu Lee



Contents

Introduction. Transprecision Computing Concepts (15 mins)

Part A. Transprecision techniques for Linear Solver (30 mins)

5 mins break;

Part B. Transprecision techniques for Kernel Recursive Least Squares (30 mins)

10 mins QnA

Contents

Beside **transprecision skills**, we will also learn **machine learning** and **linear algebra** such as:

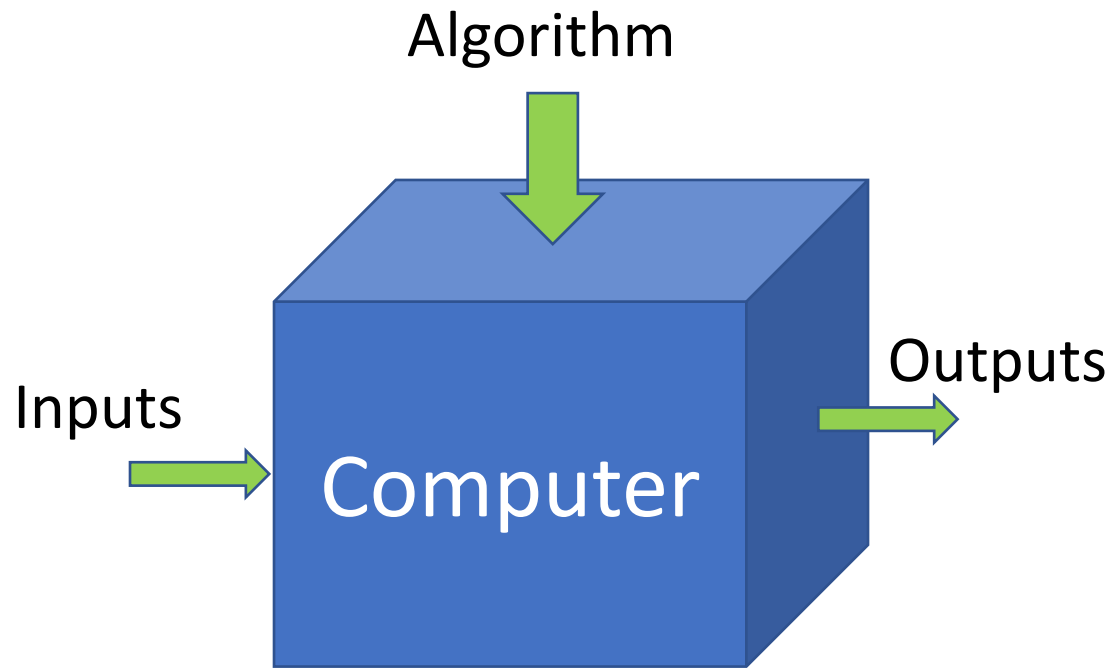
- Accuracy and Precision
- Gaussian Elimination (LU Decomposition)
- LU direct Solver
- Iterative Refinement
- Condition Number
- Global Positioning System
- Kernel Method Machine Learning

Transprecision Concept



Computation Requirements

Precision vs Accuracy?

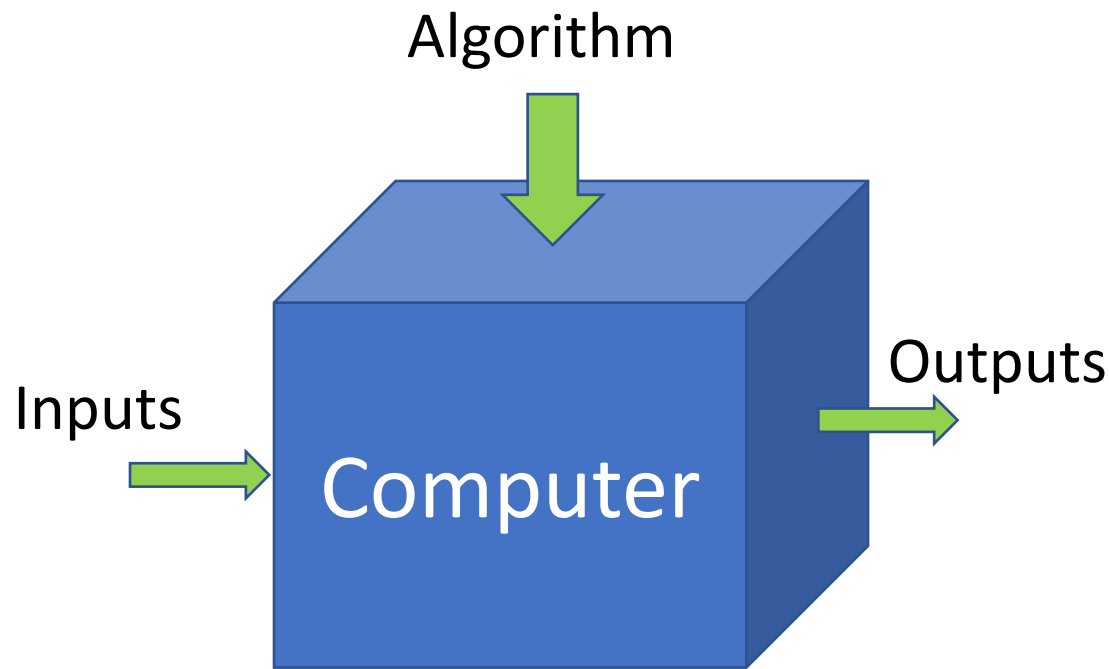


Computation Requirements

Precision vs Accuracy:

Accuracy: the absolute or relative error of an approximate quantity.

Precision: the accuracy with which the basic arithmetic operations $+$, $-$, $*$, $/$ are performed. - "Accuracy and Stability of Numerical Algorithms", NJ Higham -

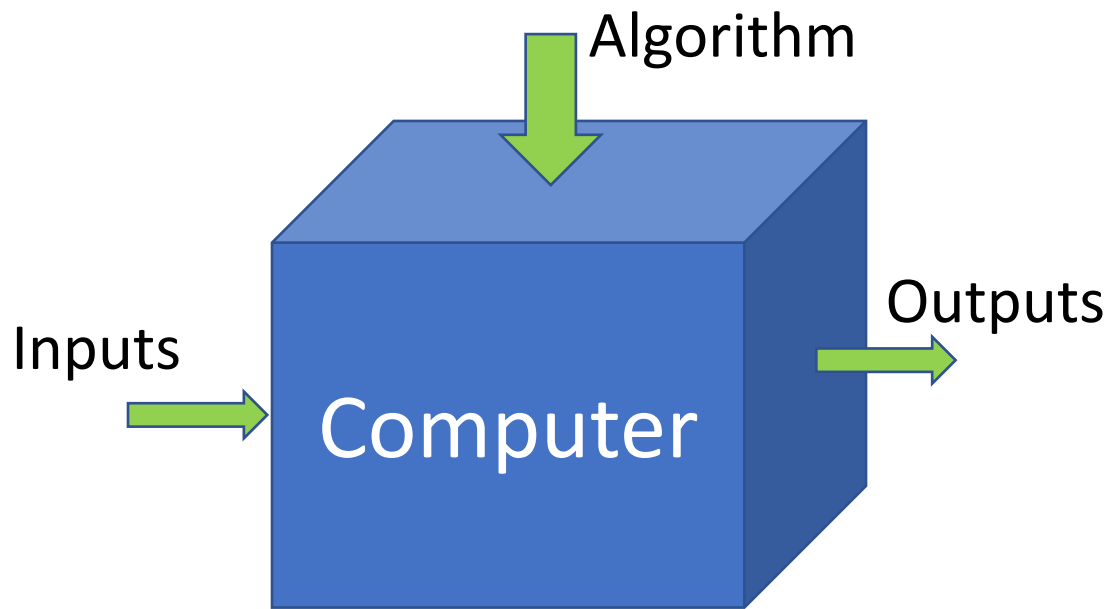


Computation Requirements

Precision vs Accuracy:

Accuracy: the absolute or relative error of an approximate quantity.

Precision: Every basic operation accuracy should be within one machine epsilon (IEEE 754). **Single Precision** (32bits) Machine Epsilon: 2^{-24} **Double Precision** (64bits) Machine Epsilon: 2^{-53}

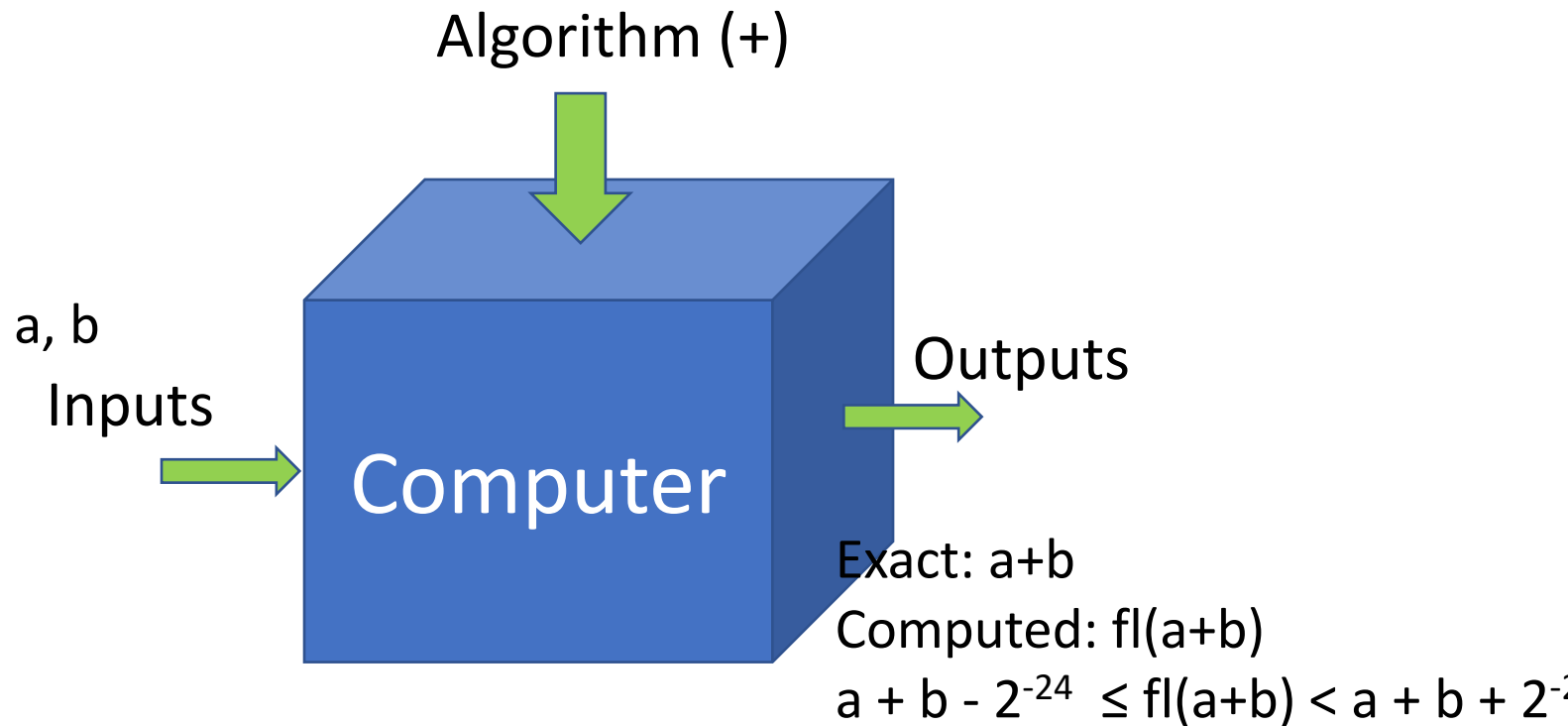


Computation Requirements

Eg. Accuracy and Precision for single precision arith Linear Solver $A \mathbf{x} = \mathbf{b}$

Accuracy: $\frac{\|x^* - x\|}{\|x^*\|}$

Precision: machine epsilon for single precision = 2^{-24} .

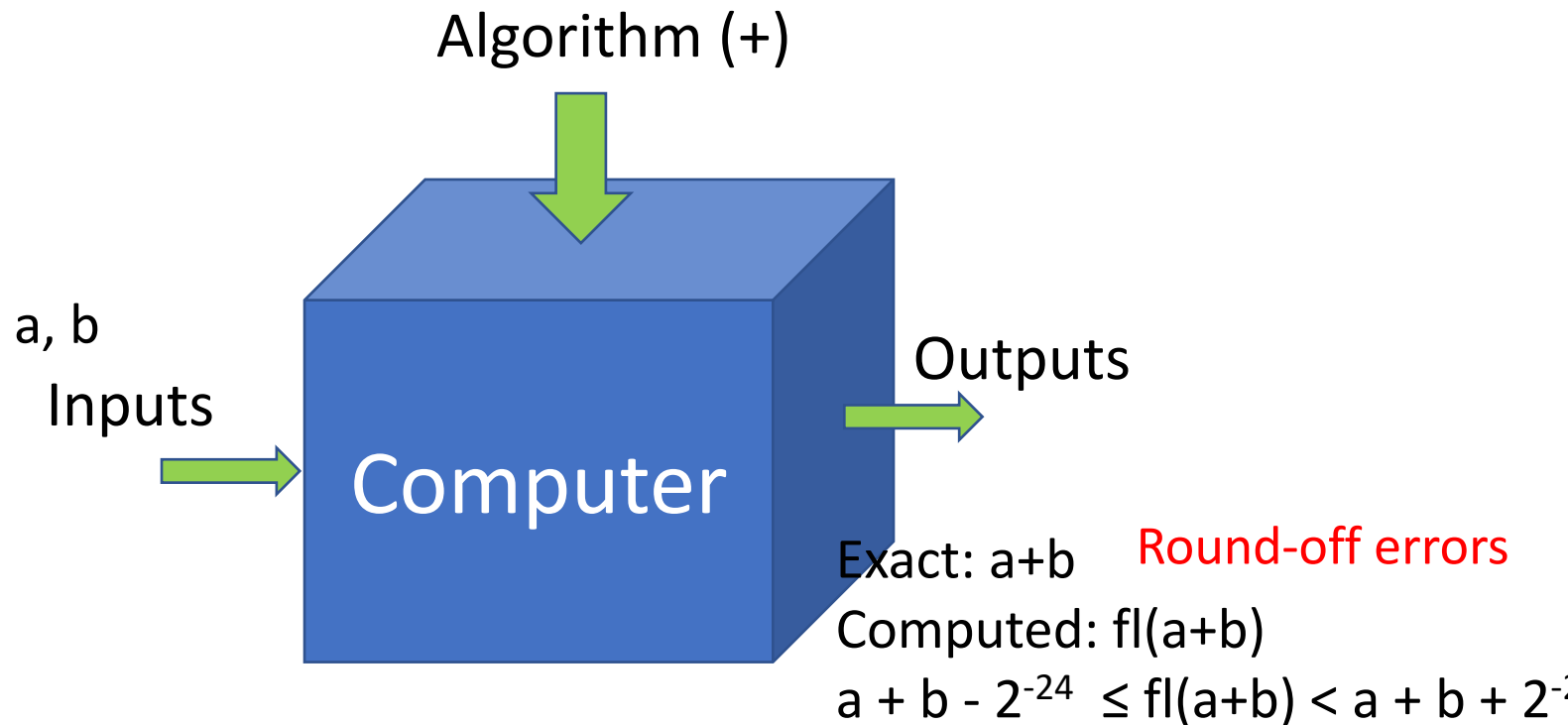


Computation Requirements

Eg. Accuracy and Precision for single precision arith Linear Solver $A \mathbf{x} = \mathbf{b}$

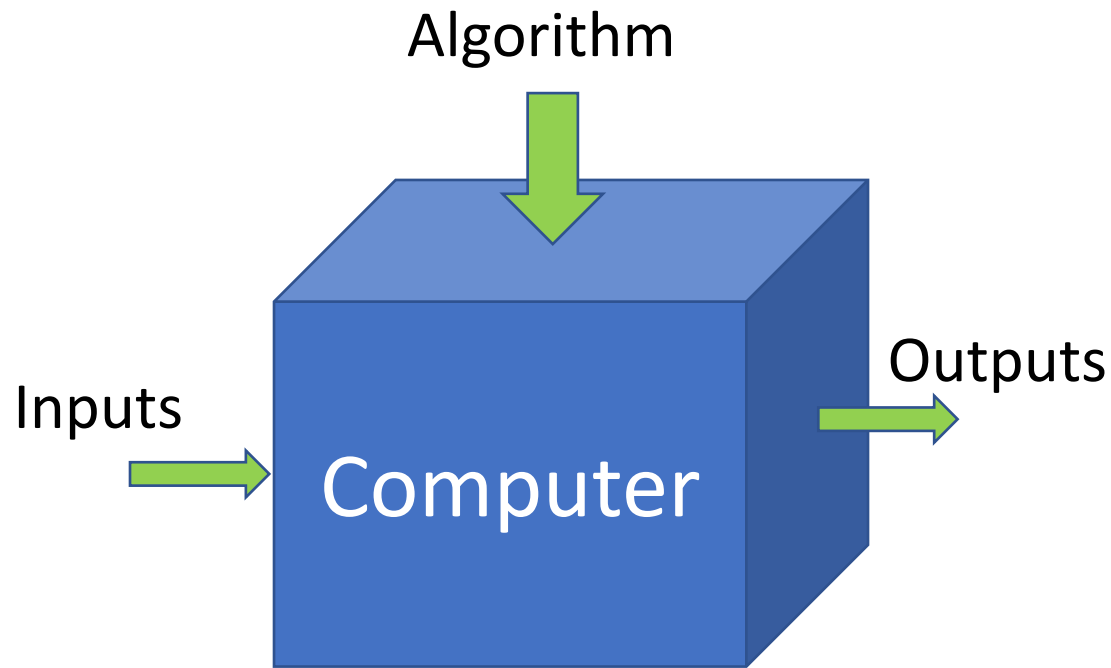
Accuracy: $\frac{\|x^* - x\|}{\|x^*\|}$

Precision: machine epsilon for single precision = 2^{-24} .



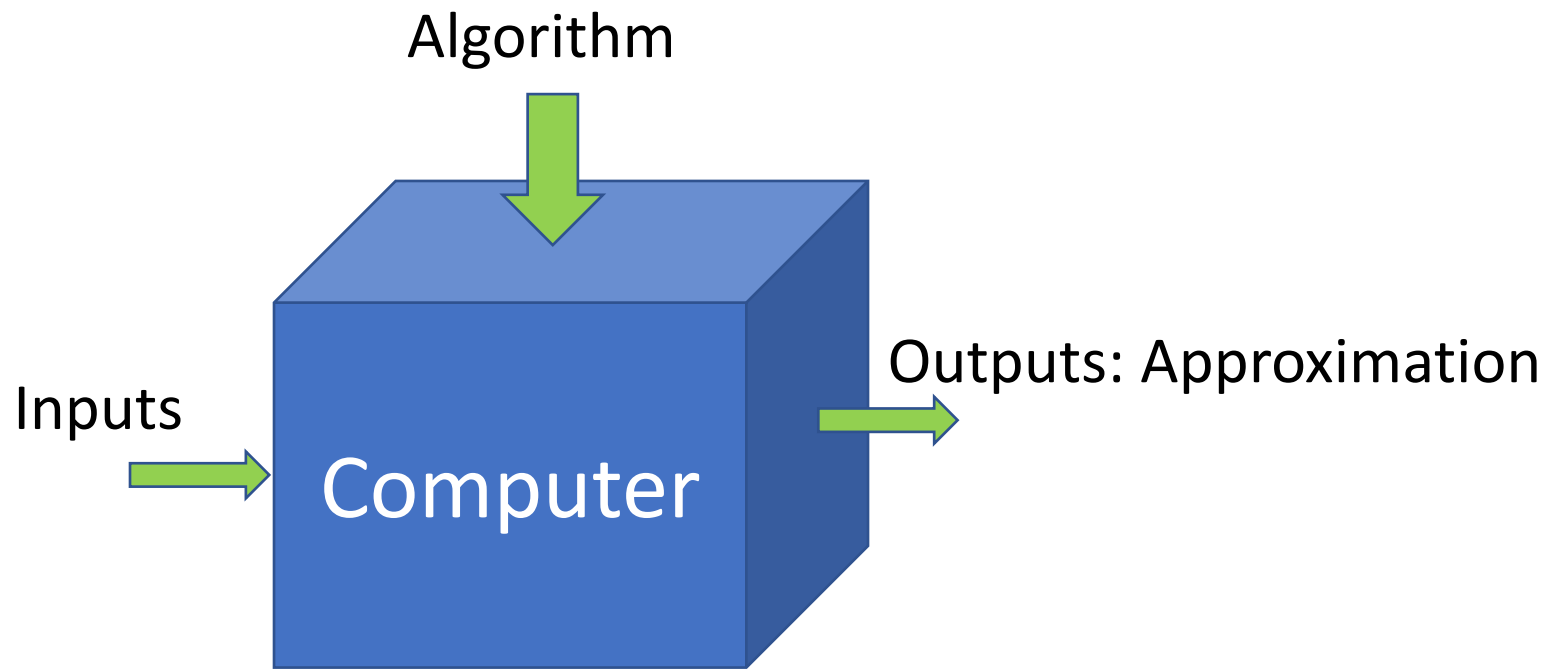
Computation Requirements

Finite Precision Arithmetic vs Exact Arithmetic ?



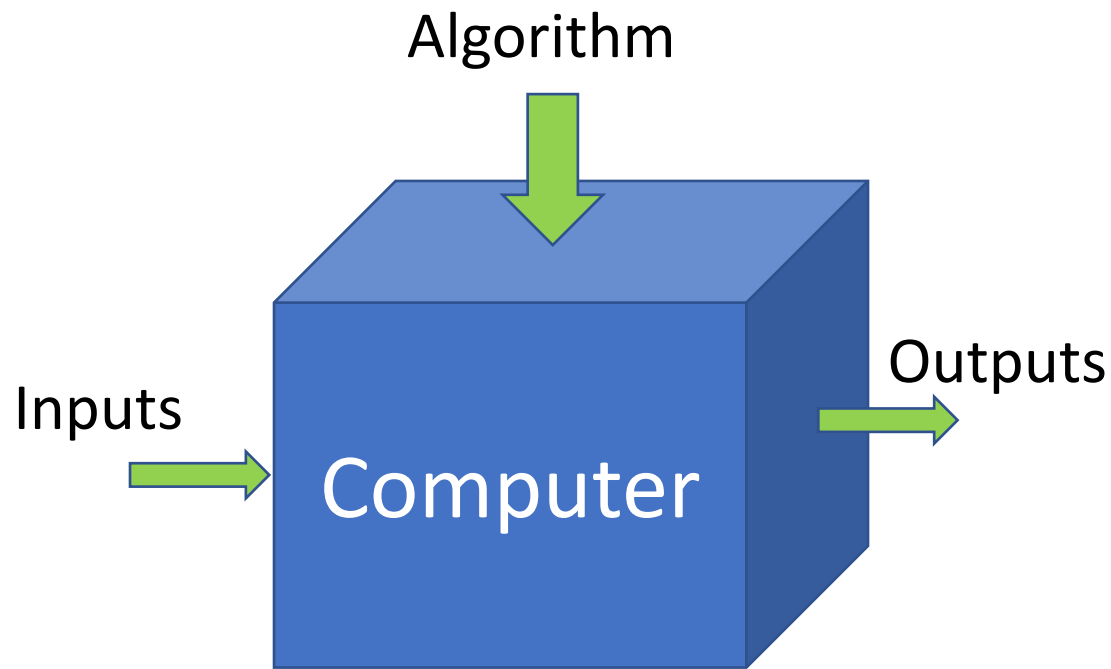
Computation Requirements

Finite Precision Arithmetic vs Exact Arithmetic ?



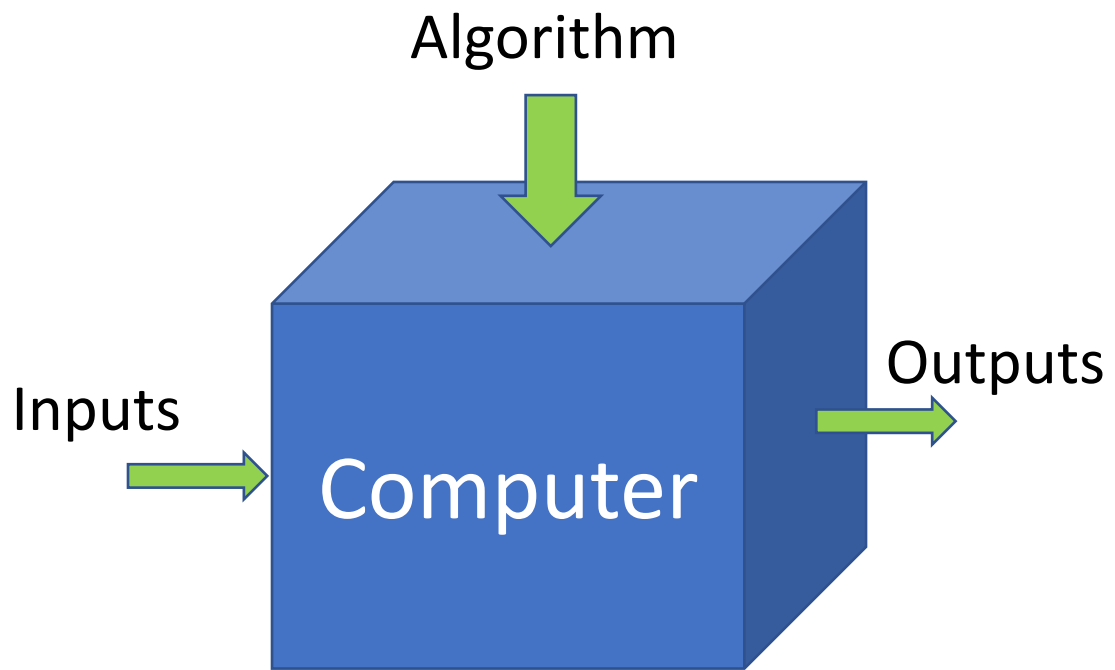
Computation Requirements

Solution Accuracy ?



Computation Requirements

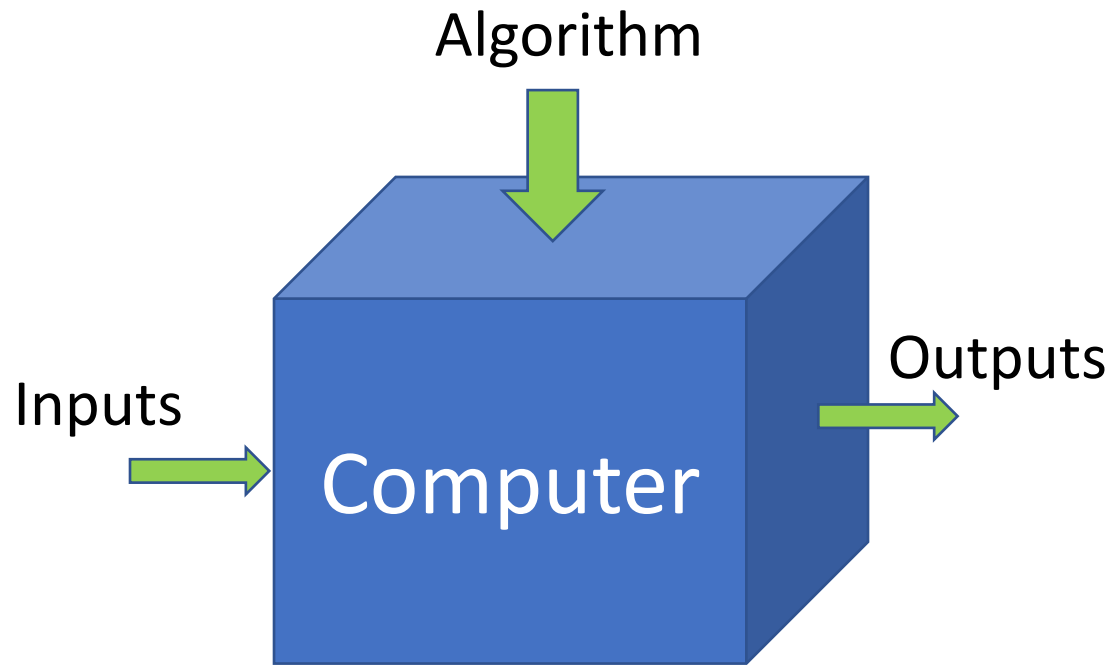
If Solution Accuracy is Good enough,
Speed and Energy?



Computation Requirements

If Solution Accuracy is Good enough,

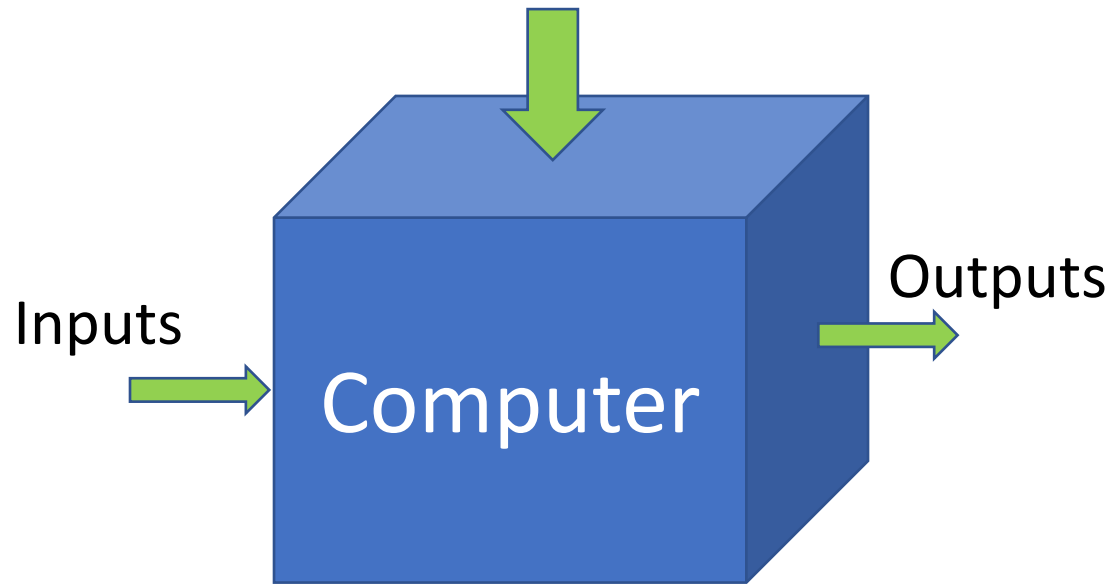
Transprecision Techniques minimizes Energy and Execution time without losing Solution Accuracy!



Computation Requirements

How?

Algorithm (Linear Solver, Non-Linear Regression)

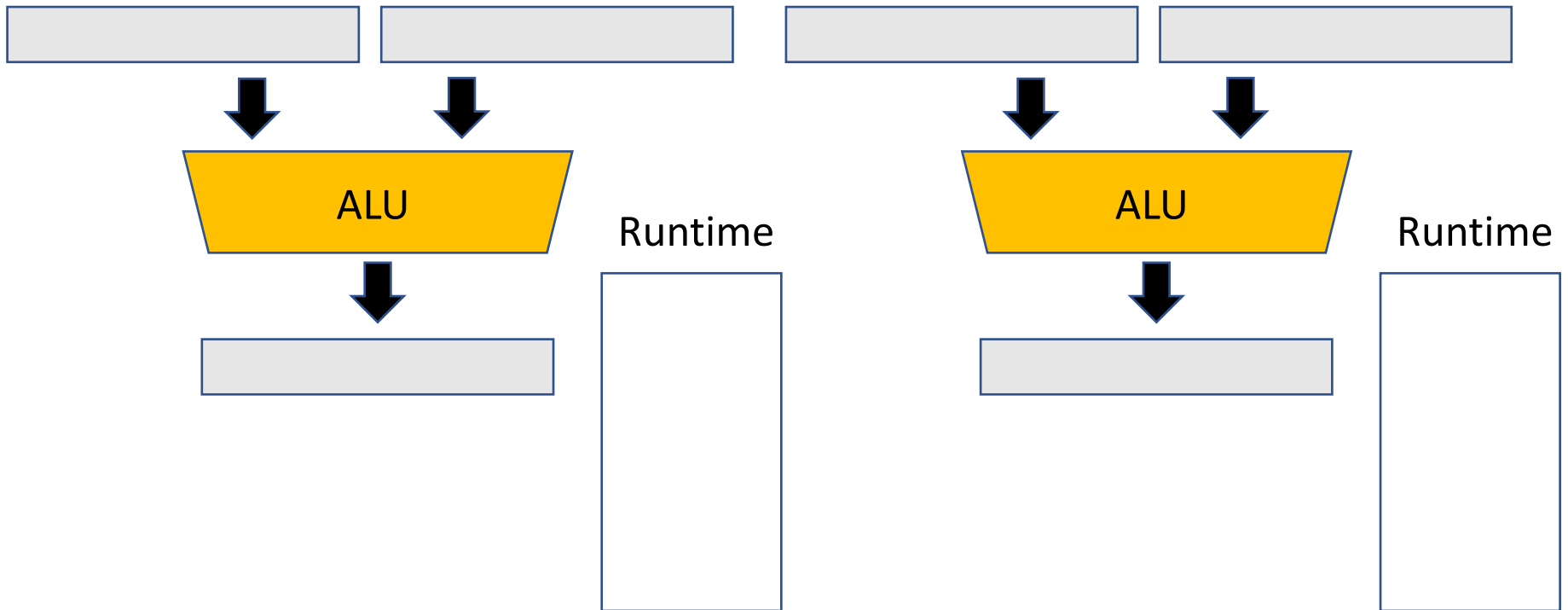


Transprecision Computing Simple Demo

Execution time for 7 operations on CPU/GPU with double precision data

Uni-precision Computing

Transprecision Computing

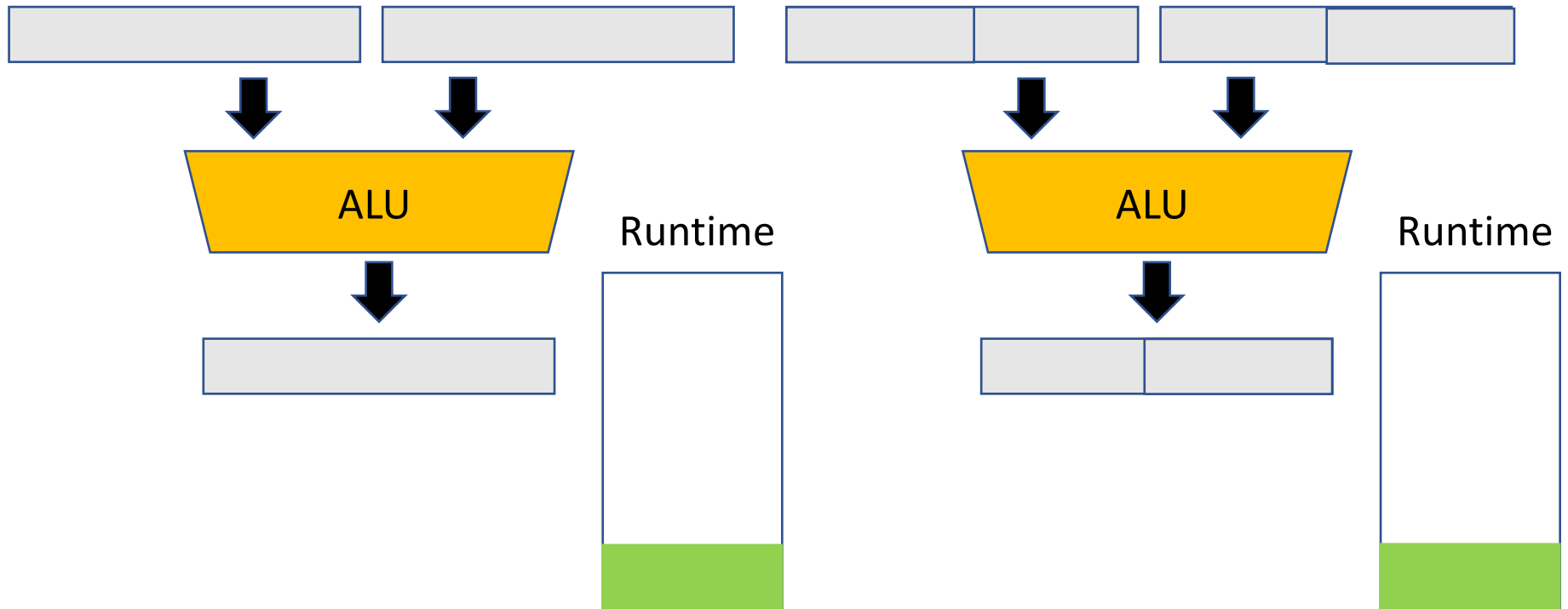


Transprecision Computing Simple Demo

Execution time for 7 operations on CPU/GPU with double precision data

Uni-precision Computing

Transprecision Computing

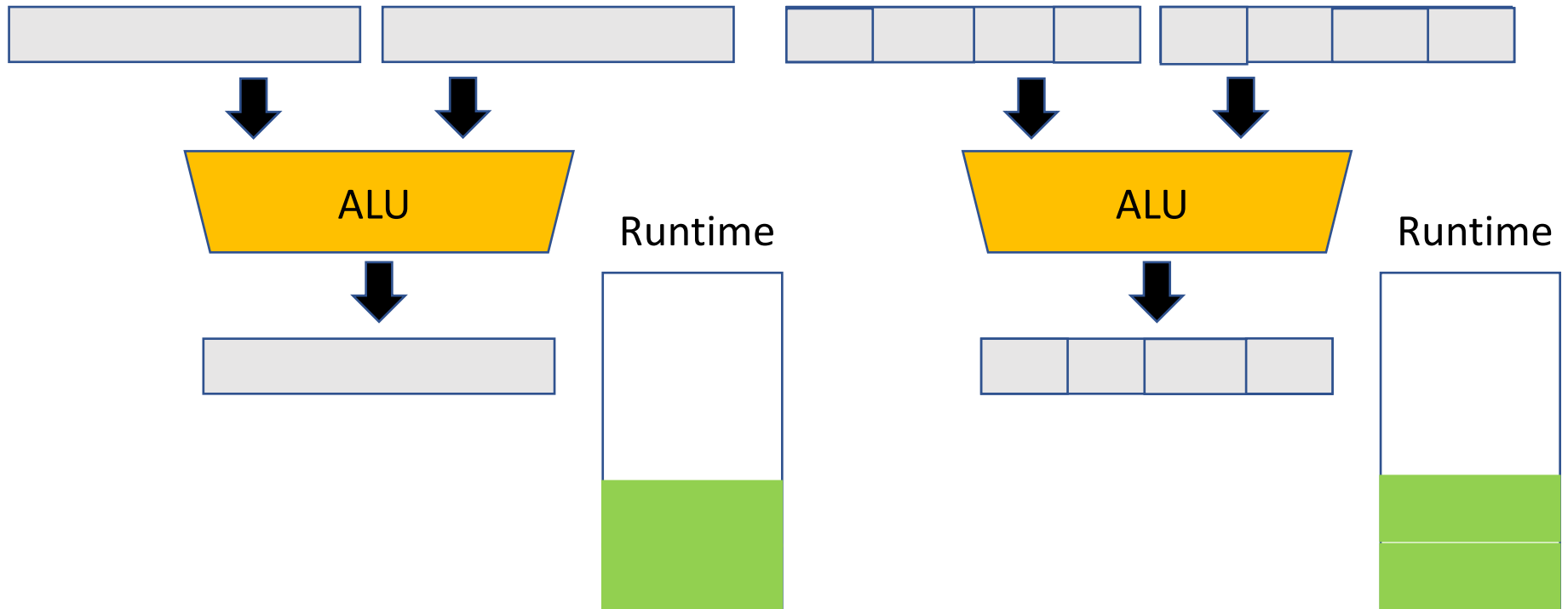


Transprecision Computing Simple Demo

Execution time for 7 operations on CPU/GPU with double precision data

Uni-precision Computing

Transprecision Computing

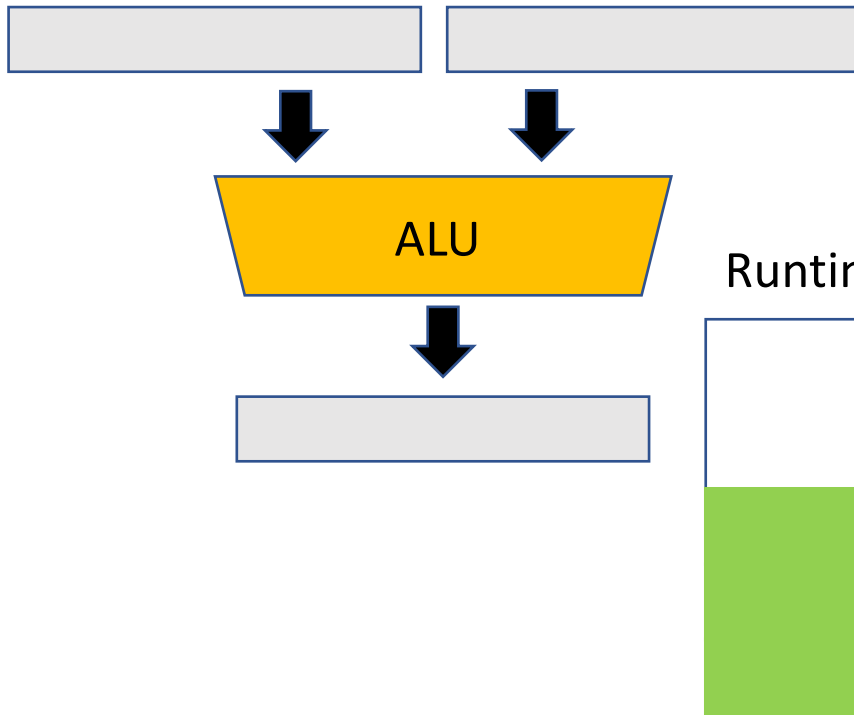


Transprecision Computing Simple Demo

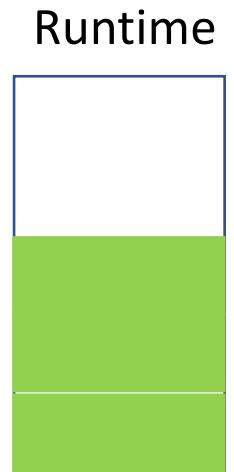
Execution time for 7 operations on CPU/GPU with double precision data

Uni-precision Computing

Transprecision Computing



DONE!

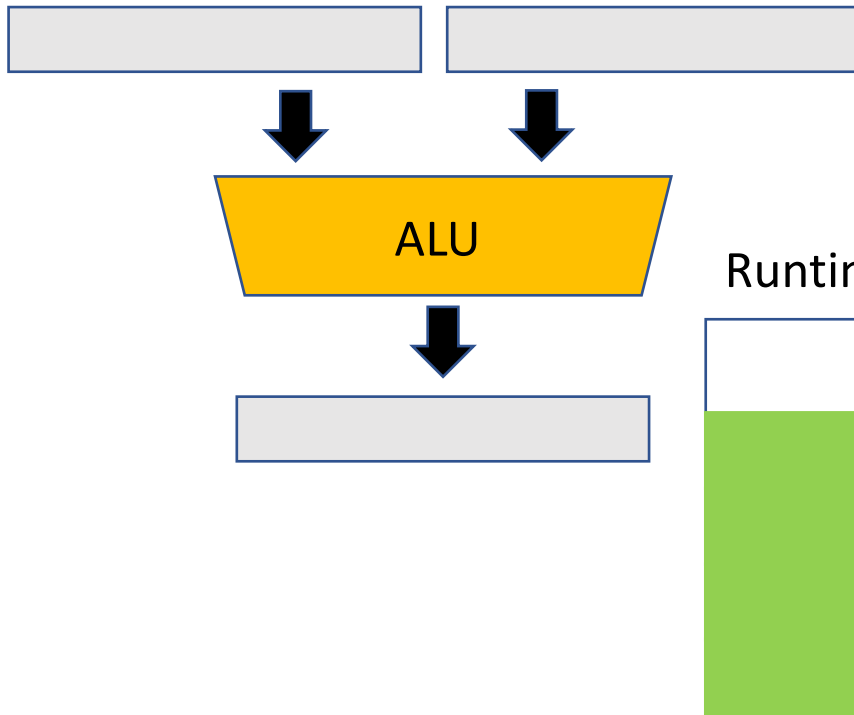


Transprecision Computing Simple Demo

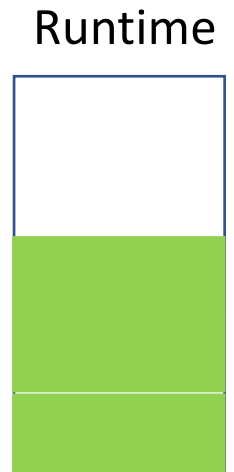
Execution time for 7 operations on CPU/GPU with double precision data

Uni-precision Computing

Transprecision Computing



DONE!

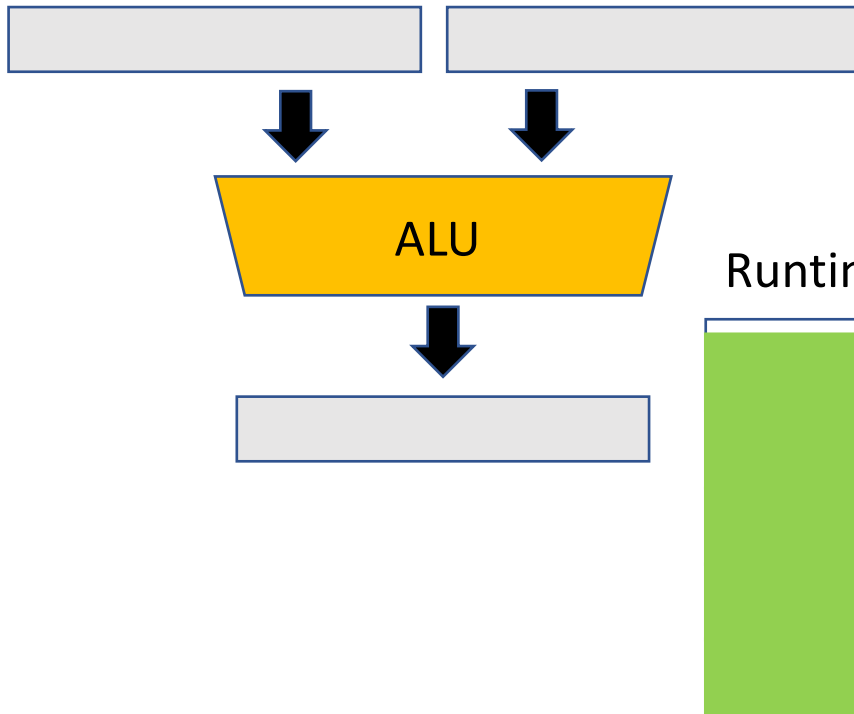


Transprecision Computing Simple Demo

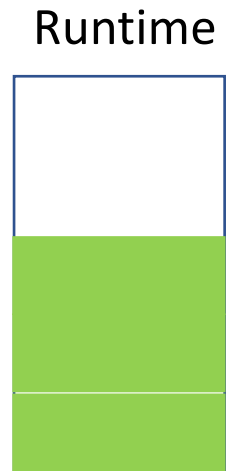
Execution time for 7 operations on CPU/GPU with double precision data

Uni-precision Computing

Transprecision Computing



DONE!

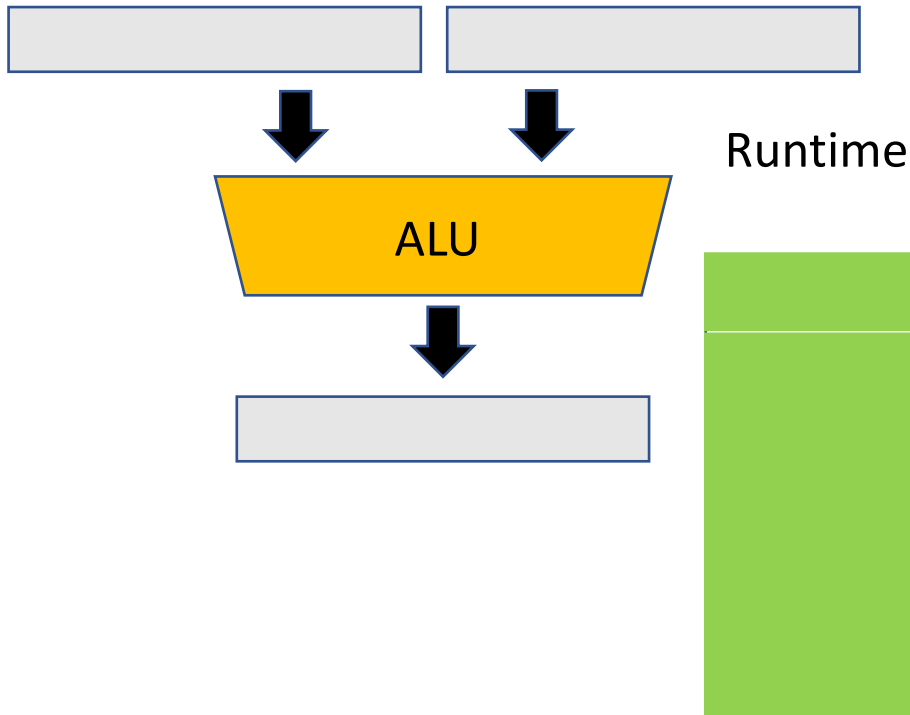


Transprecision Computing Simple Demo

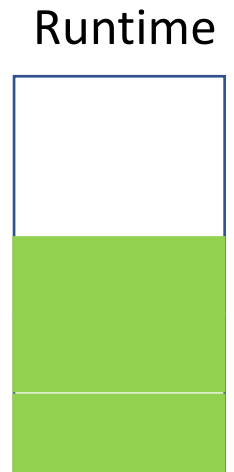
Execution time for 7 operations on CPU/GPU with double precision data

Uni-precision Computing

Transprecision Computing



DONE!



Transprecision Computing Simple Demo

Execution time for 7 operations on CPU/GPU with double precision data

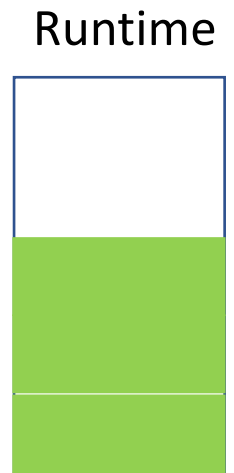
Uni-precision Computing

DONE!



Transprecision Computing

DONE!

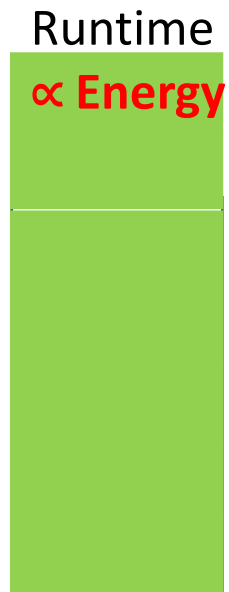


Transprecision Computing Simple Demo

Execution time for 7 operations on CPU/GPU with double precision data

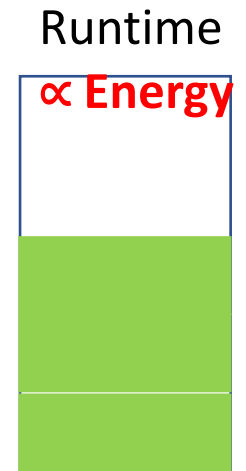
Uni-precision Computing

DONE!



Transprecision Computing

DONE!



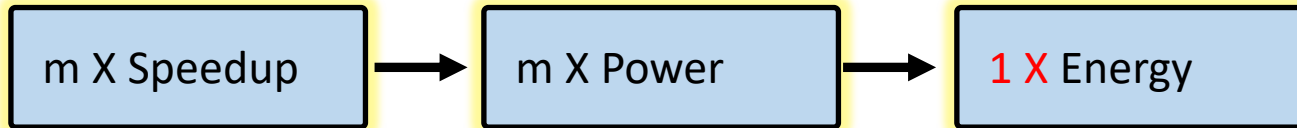
Transprecision Techniques seek

minimised precision arithmetic without losing solution accuracy

Transprecision Computing vs Parallel Computing

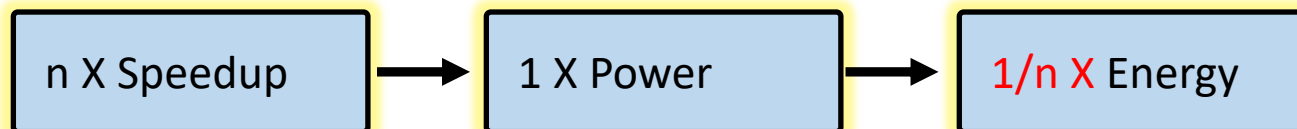
Reduced runtime without increasing number of cores
=> Direct Energy Savings!

Parallel Computing with $m \times$ Cores



Need some techniques for energy saving?

Transprecision Computing without increasing cores



Energy Savings!

Transprecision Impact on GPUs/FPGAs

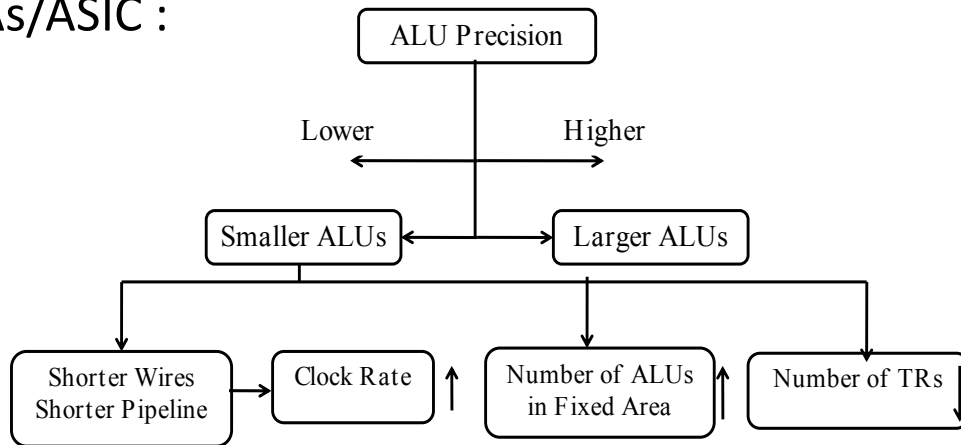
NVidia Pascal GPU (P100) with NVLink:

- Half precision: 21.2 TeraFlops
- Single precision: 10.6 TeraFlops
- Double precision: 5.3 TeraFlops

NVidia Volta GPU (V100) with NVLink:

- Half precision (Tensor Core): 125 TeraFlops
- Single precision: 15.7 TeraFlops
- Double precision: 7.8 TeraFlops

FPGAs/ASIC :



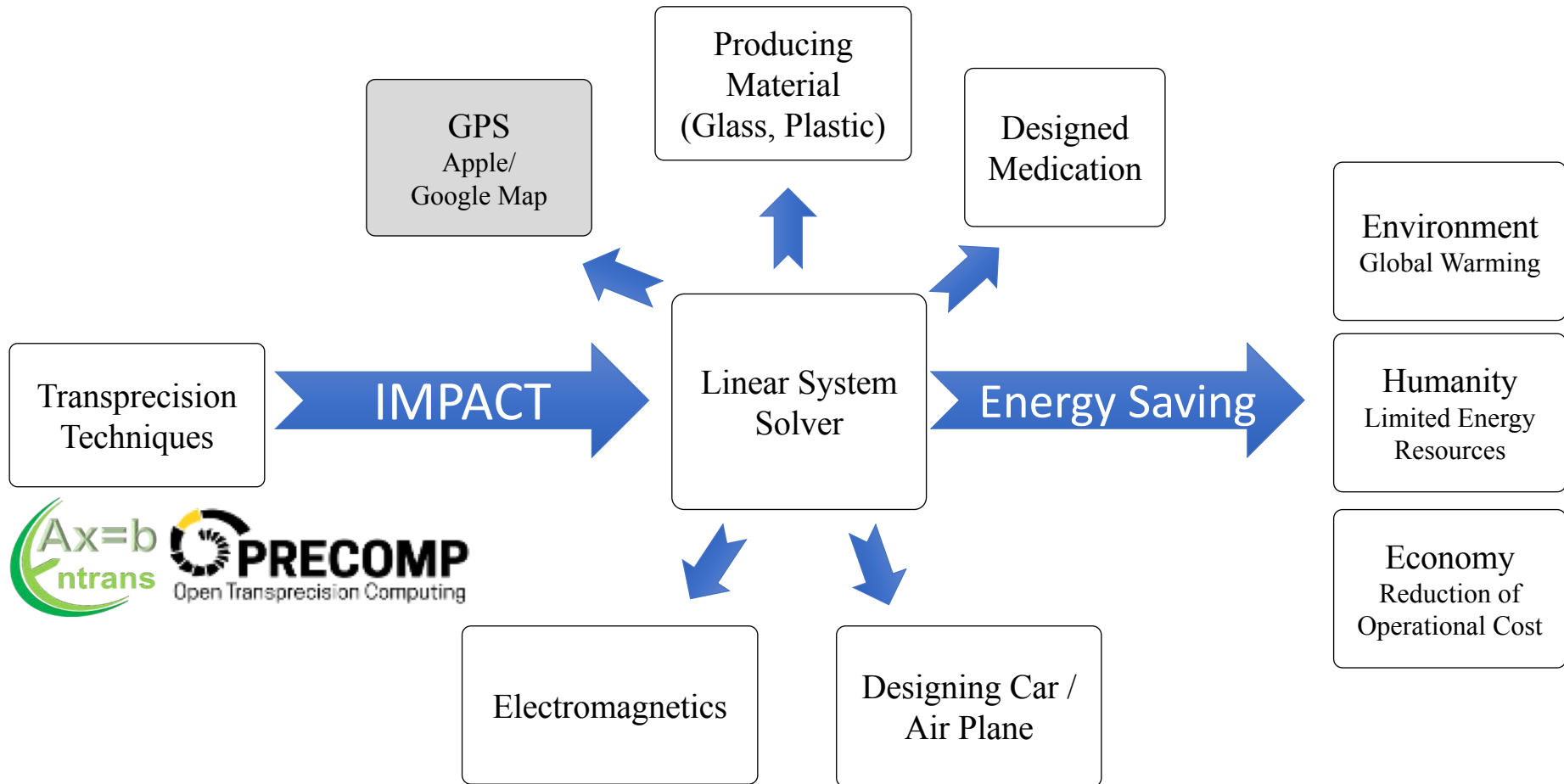
Size of Adder: Linear increase with precision

Size of Multiplier: Quadratic increase with precision

Part A. Linear Solver



Linear Solver Applications



Global Positioning System (GPS)

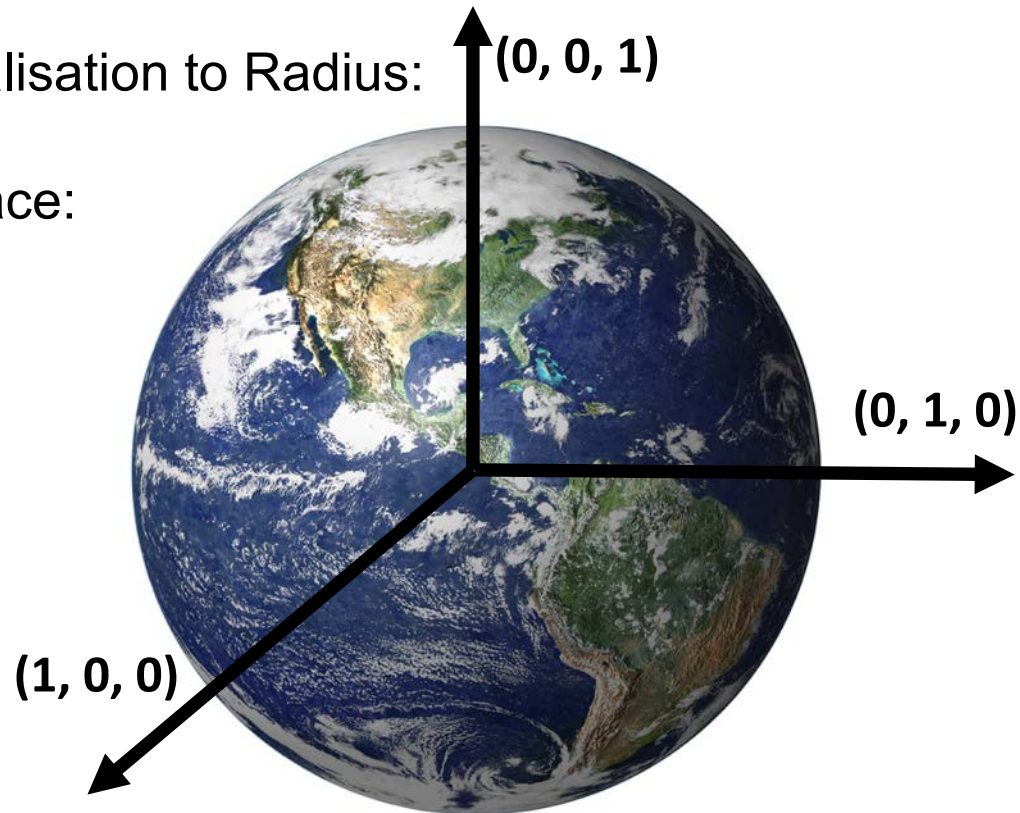
Source (reference):

<http://www.math.uri.edu/~jbaglama/classes/2011-2012/spring/mth215/project1.pdf>

Positioning Normalisation to Radius:

Any point on Surface:

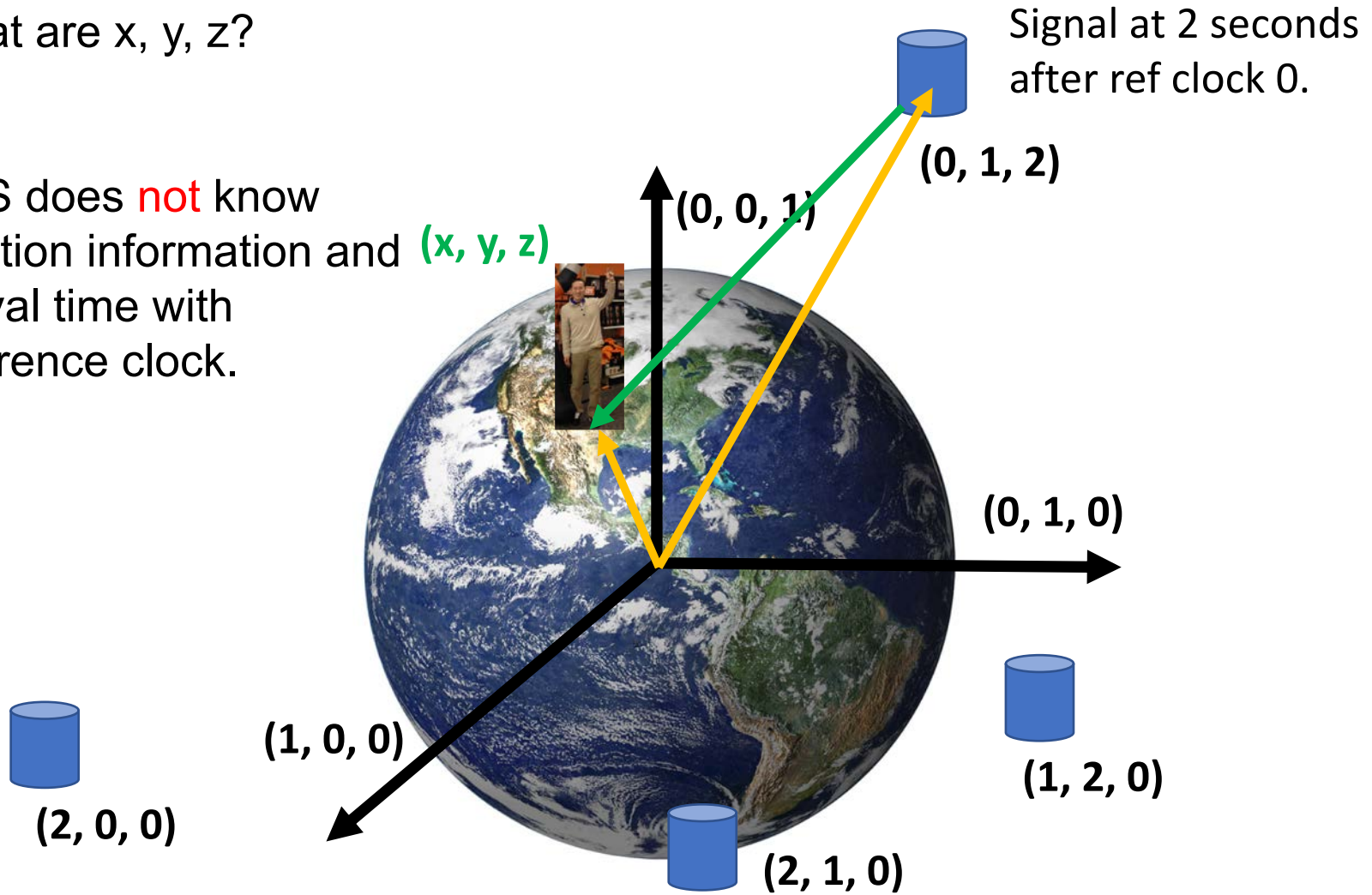
$$x^2 + y^2 + z^2 = 1$$



Global Positioning System (GPS)

What are x, y, z ?

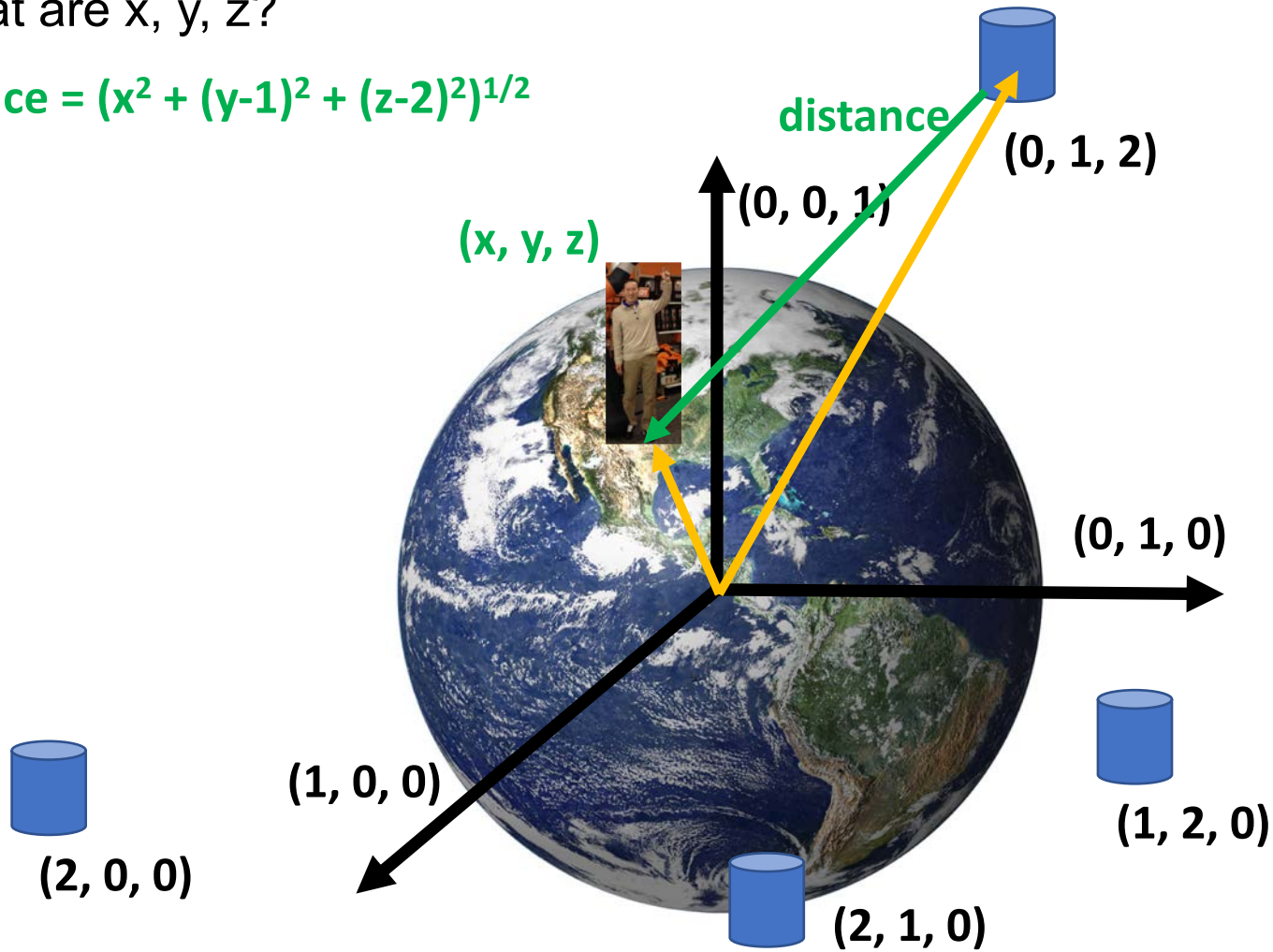
GPS does **not** know location information and (x, y, z) arrival time with reference clock.



Global Positioning System (GPS)

What are x, y, z?

$$\text{distance} = (x^2 + (y-1)^2 + (z-2)^2)^{1/2}$$



Global Positioning System (GPS)

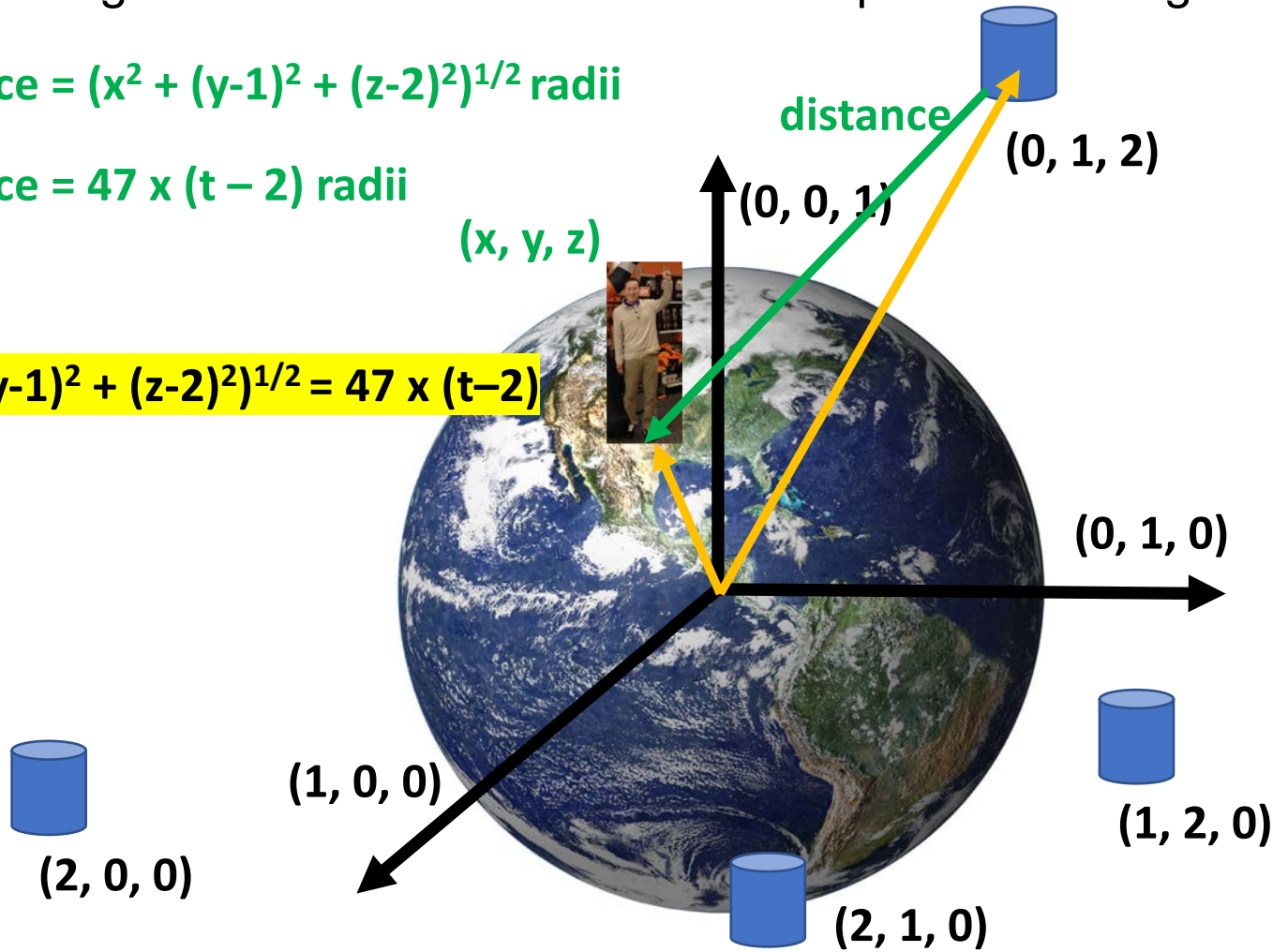
Assuming $c = 47$ radii/sec and the satellite provides a single in $t=2$ seconds.

$$\text{distance} = (x^2 + (y-1)^2 + (z-2)^2)^{1/2} \text{ radii}$$

$$\text{distance} = 47 \times (t - 2) \text{ radii}$$

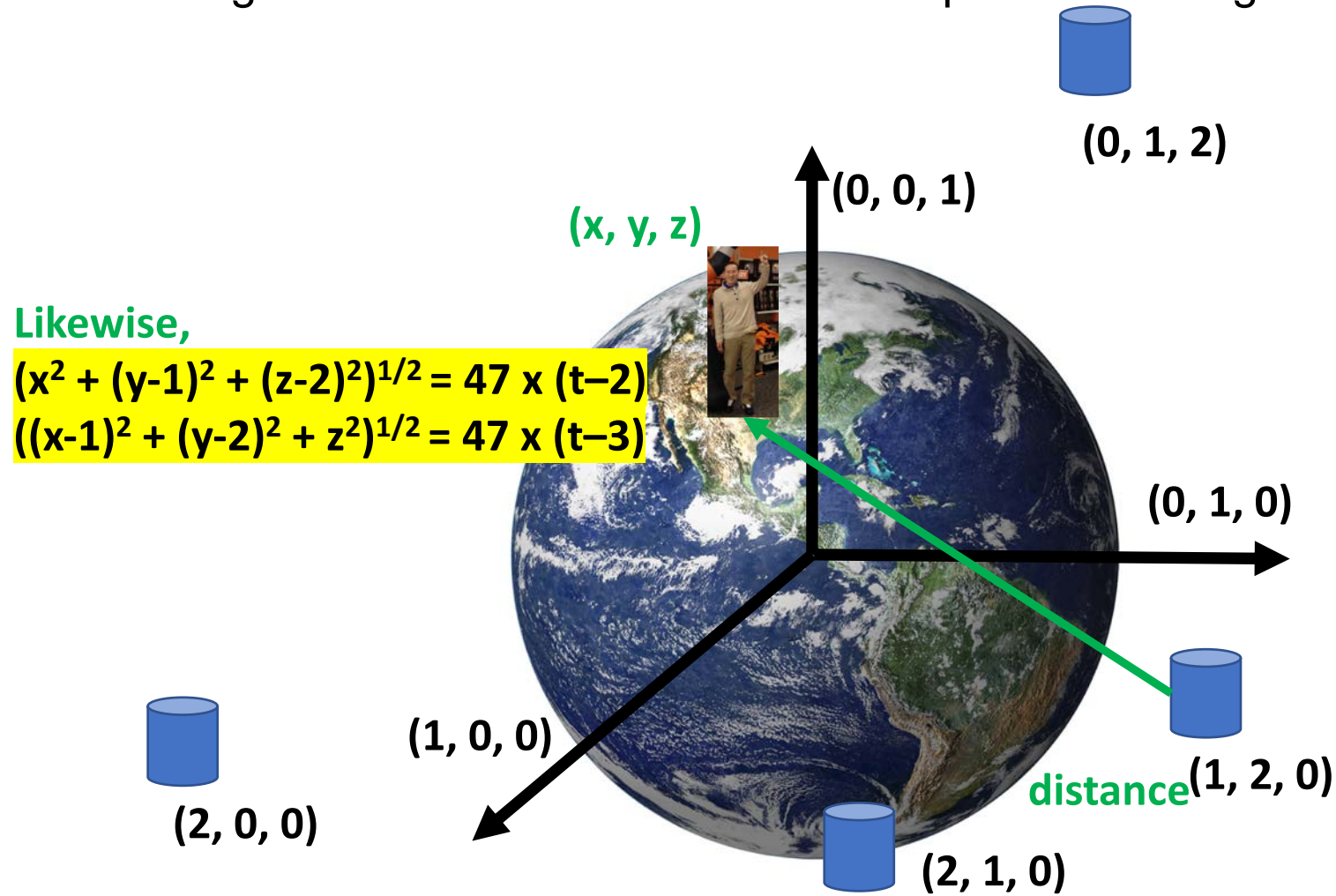
(x, y, z)

$$(x^2 + (y-1)^2 + (z-2)^2)^{1/2} = 47 \times (t-2)$$



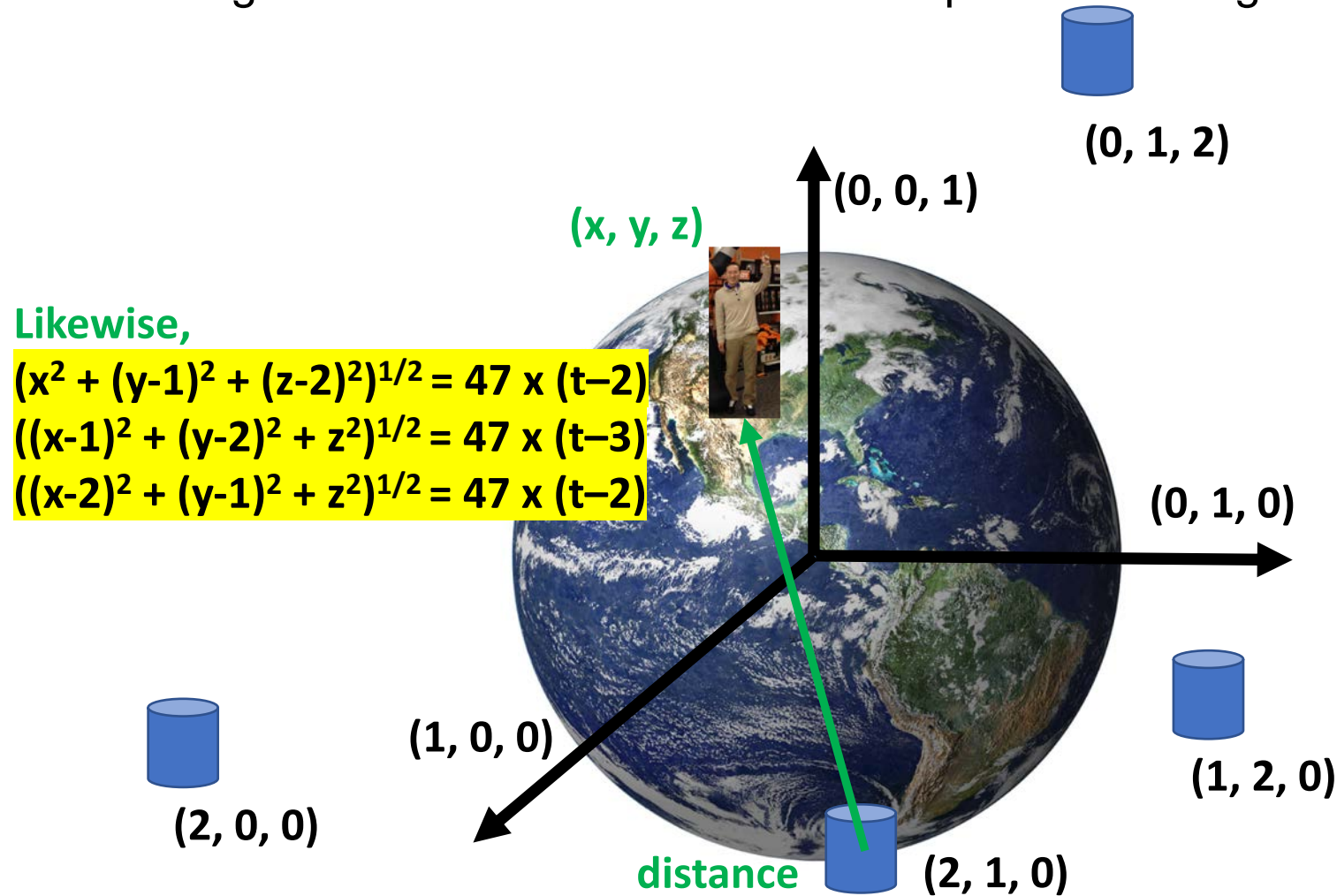
Global Positioning System (GPS)

Assuming $c = 47$ radii/sec and the satellite provides a single in $t=2$ seconds.



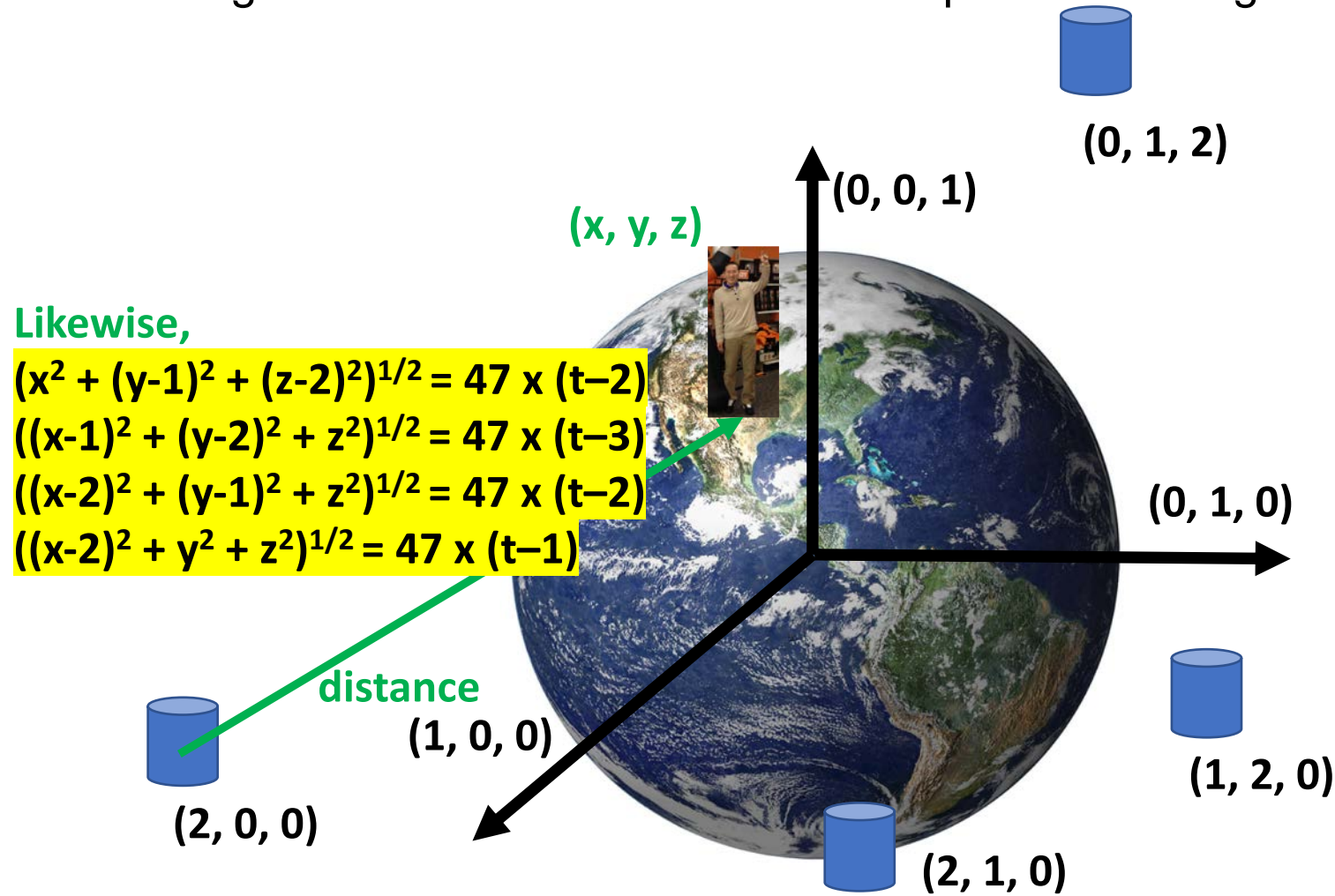
Global Positioning System (GPS)

Assuming $c = 47$ radii/sec and the satellite provides a single in $t=2$ seconds.



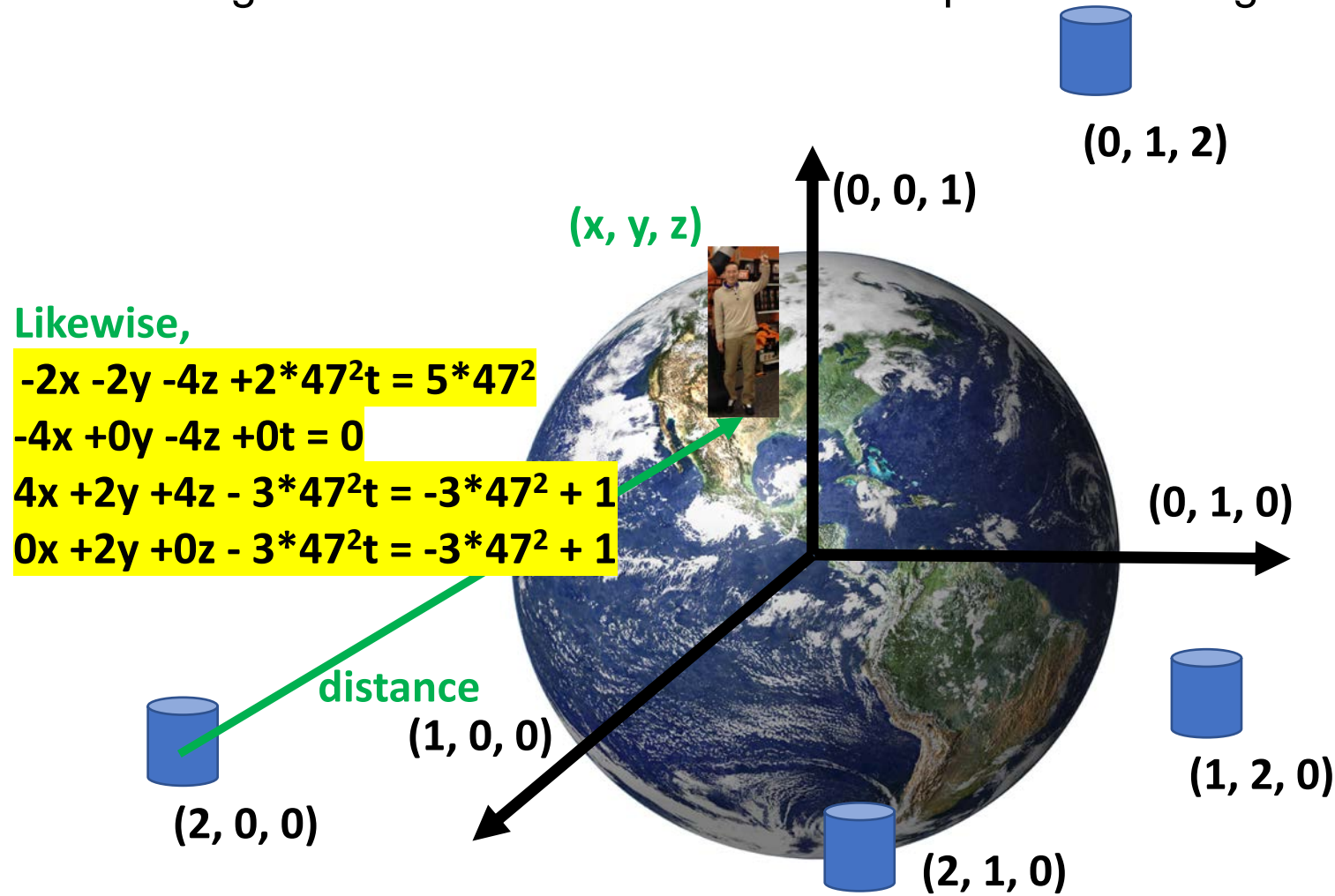
Global Positioning System (GPS)

Assuming $c = 47$ radii/sec and the satellite provides a single in $t=2$ seconds.



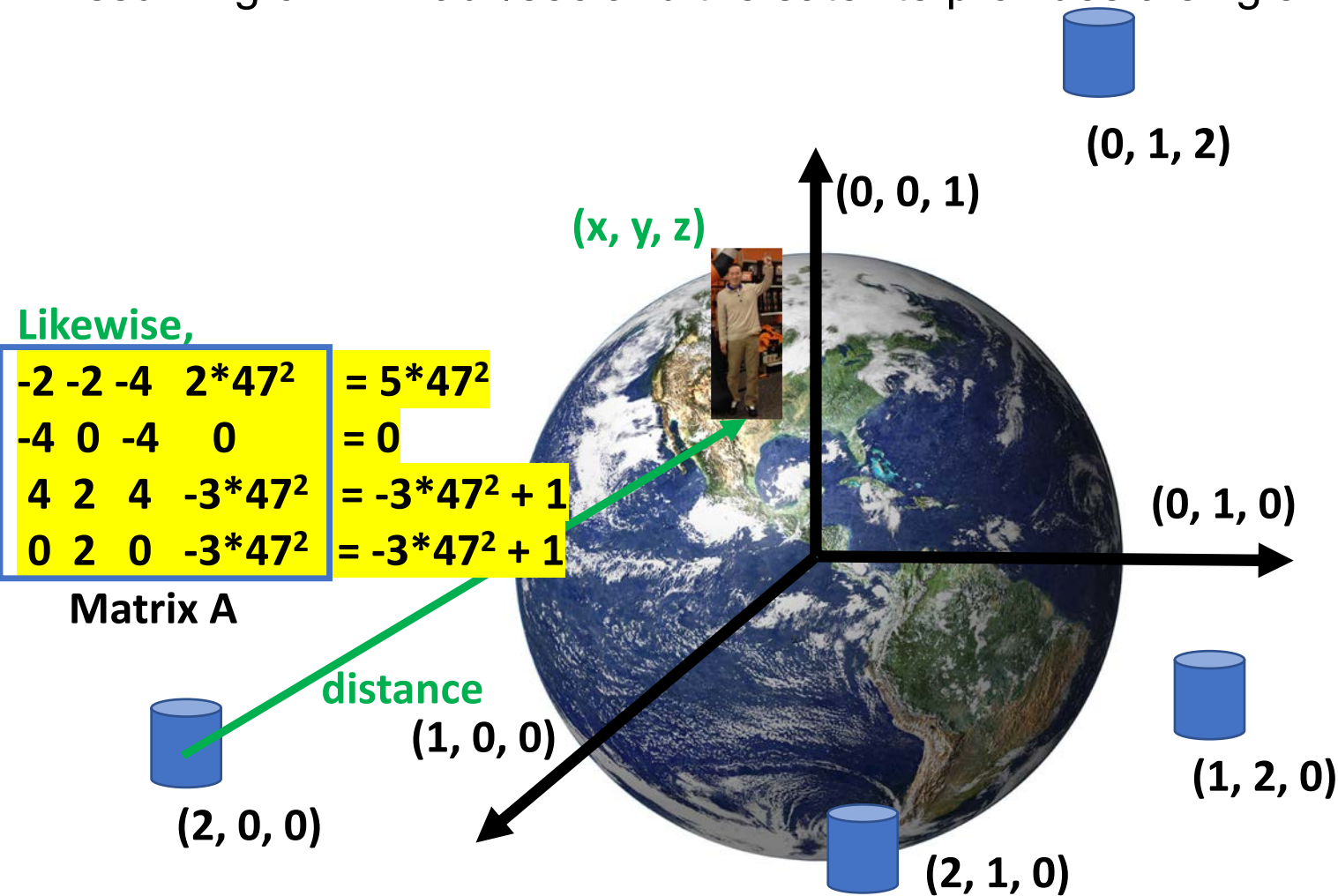
Global Positioning System (GPS)

Assuming $c = 47$ radii/sec and the satellite provides a single in $t=2$ seconds.



Global Positioning System (GPS)

Assuming $c = 47$ radii/sec and the satellite provides a single in $t=2$ seconds.



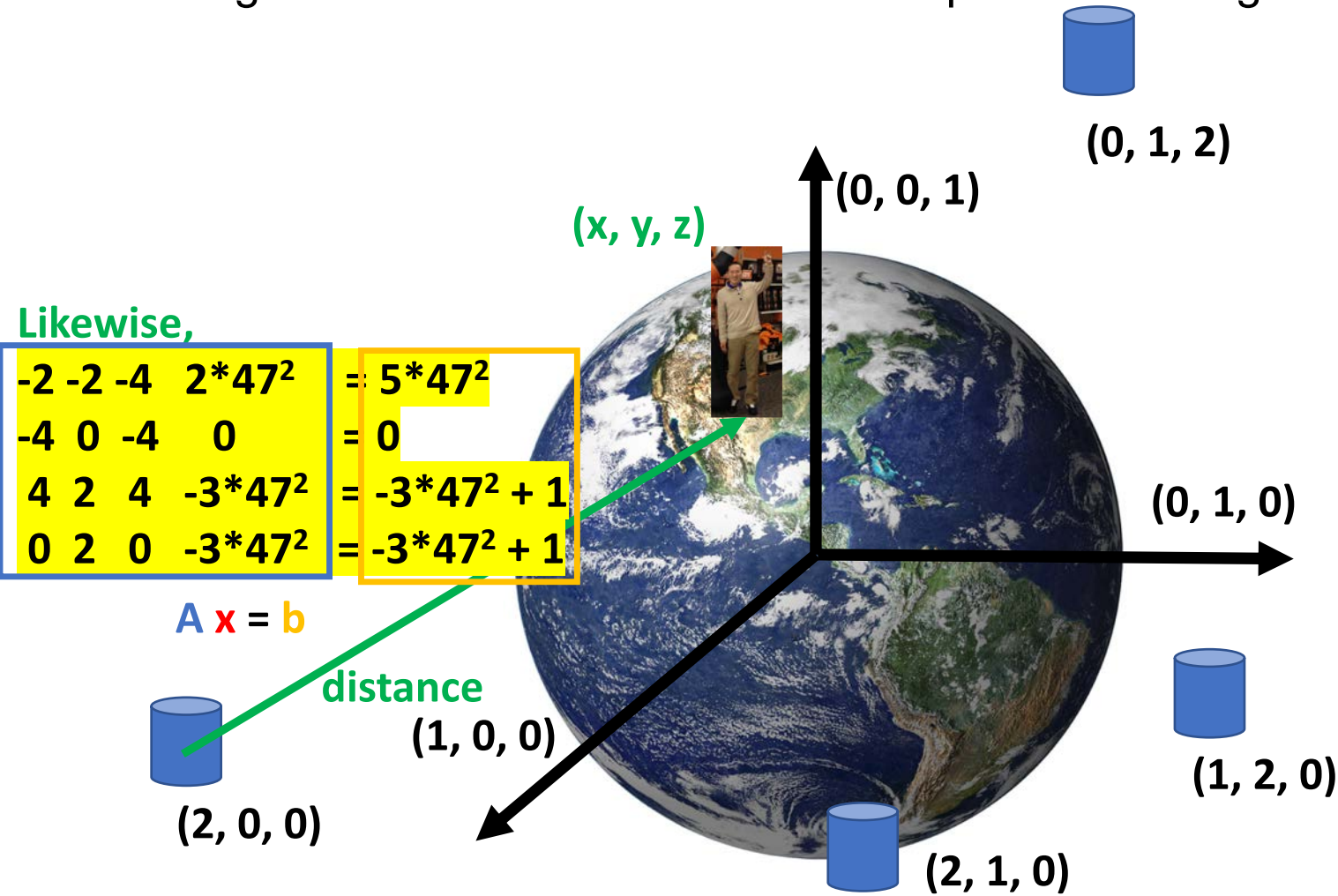
Likewise,

-2	-2	-4	$2 \cdot 47^2$	$= 5 \cdot 47^2$
-4	0	-4	0	$= 0$
4	2	4	$-3 \cdot 47^2$	$= -3 \cdot 47^2 + 1$
0	2	0	$-3 \cdot 47^2$	$= -3 \cdot 47^2 + 1$

Matrix A

Global Positioning System (GPS)

Assuming $c = 47$ radii/sec and the satellite provides a single in $t=2$ seconds.



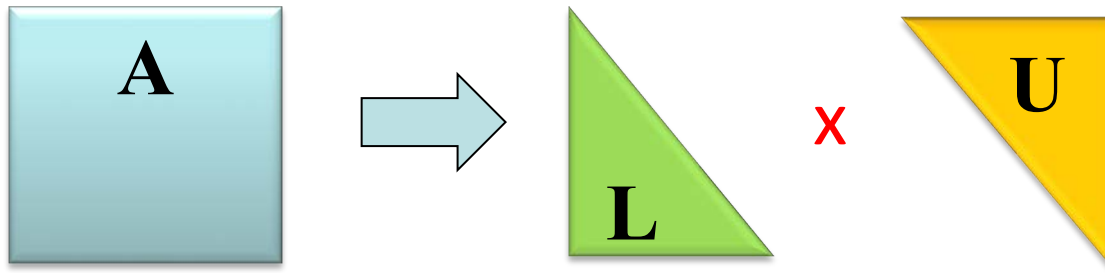
Linear Solver Algorithm

Find x in $Ax = b$!

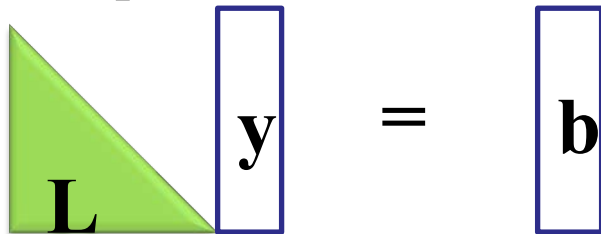
Direct Solver ($A\mathbf{x} = \mathbf{b}$)

LU decomposition

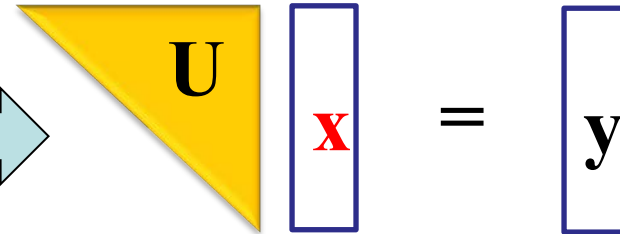
Matrix Decomposition: $\frac{2}{3} n^3$ Ops



n^2 Ops (Forward Substitution)



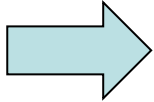
n^2 Ops (Back Substitution)



LU Factorisation for 3 x 3 Matrix (i.e., Gaussian Elimination)

Gaussian Elimination eliminates a_{21}, a_{31}, a_{32}

2	3	1
5	1	3
0	2	4

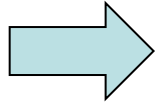


1	0	0
l_{21}	1	0
l_{31}	l_{32}	1

u_{11}	u_{12}	u_{13}
0	u_{22}	u_{23}
0	0	u_{33}

At 2nd Row, multiply $5/2$ to 1st Row and Subtract ; Store $l_{21} \leftarrow 5/2$

2	3	1
5	1	3
0	2	4



2	3	1
0	$-13/2$	$1/2$
0	2	4

LU Factorisation for 3 x 3 Matrix (i.e., Gaussian Elimination)

Gaussian Elimination eliminates a_{21}, a_{31}, a_{32}

$$\begin{array}{|c|c|c|} \hline 2 & 3 & 1 \\ \hline 5 & 1 & 3 \\ \hline 0 & 2 & 4 \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline 5/2 & 1 & 0 \\ \hline l_{31} & l_{32} & 1 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline u_{11} & u_{12} & u_{13} \\ \hline 0 & u_{22} & u_{23} \\ \hline 0 & 0 & u_{33} \\ \hline \end{array}$$

At 3rd Row, multiply 0 to 1st Row and Subtract ; Store $l_{31} \leftarrow 0$

$$\begin{array}{|c|c|c|} \hline 2 & 3 & 1 \\ \hline 0 & -13/2 & 1/2 \\ \hline 0 & 2 & 4 \\ \hline \end{array} \rightarrow \begin{array}{|c|c|c|} \hline 2 & 3 & 1 \\ \hline 0 & -13/2 & 1/2 \\ \hline 0 & 2 & 4 \\ \hline \end{array}$$

LU Factorisation for 3 x 3 Matrix (i.e., Gaussian Elimination)

Gaussian Elimination eliminates a_{21}, a_{31}, a_{32}

2	3	1	→	1	0	0	u_{11}	u_{12}	u_{13}
5	1	3		$5/2$	1	0	0	u_{22}	u_{23}
0	2	4		0	l_{32}	1	0	0	u_{33}

At 3rd Row, multiply $-4/13$ to 2nd Row and Subtract ; Store $l_{32} \leftarrow -4/13$

2	3	1	→	2	3	1
0	$-13/2$	$1/2$		0	$-13/2$	$1/2$
0	2	4		0	0	$54/13$

LU Factorisation for 3 x 3 Matrix (i.e., Gaussian Elimination)

Gaussian Elimination eliminates a_{21} , a_{31} , a_{32}

2	3	1	→	1	0	0	2	3	1
5	1	3		5/2	1	0	0	-13/2	1/2
0	2	4		0	-4/13	1	0	0	54/13

At 3rd Row, multiply $-4/13$ to 2st Row and Subtract ; Store $l_{32} \leftarrow -4/13$

2	3	1	→	2	3	1
0	-13/2	1/2		0	-13/2	1/2
0	2	4		0	0	54/13

Direct LU Solver ($A\mathbf{x} = \mathbf{b}$)

Solve :

$$\begin{bmatrix} 2 & 3 & 1 \\ 5 & 1 & 3 \\ 0 & 2 & 4 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}$$

Direct LU Solver ($A\mathbf{x} = \mathbf{b}$)

Solve :

$$\begin{bmatrix} 2 & 3 & 1 \\ 5 & 1 & 3 \\ 0 & 2 & 4 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} \Leftrightarrow$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 5/2 & 1 & 0 \\ 0 & -4/13 & 1 \end{bmatrix} \begin{bmatrix} 2 & 3 & 1 \\ 0 & -13/2 & 1/2 \\ 0 & 0 & 54/13 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}$$

Direct LU Solver ($A\mathbf{x} = \mathbf{b}$)

Solve : seek y_1, y_2, y_3 by forward substitutions

$$\begin{bmatrix} 2 & 3 & 1 \\ 5 & 1 & 3 \\ 0 & 2 & 4 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 5/2 & 1 & 0 \\ 0 & -4/13 & 1 \end{bmatrix} \begin{bmatrix} 2 & 3 & 1 \\ 0 & -13/2 & 1/2 \\ 0 & 0 & 54/13 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} \Leftrightarrow$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 5/2 & 1 & 0 \\ 0 & -4/13 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}$$

Direct LU Solver ($A\mathbf{x} = \mathbf{b}$)

Solve : seek x_1, x_2, x_3 by backward substitutions

$$\begin{bmatrix} 2 & 3 & 1 \\ 5 & 1 & 3 \\ 0 & 2 & 4 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 5/2 & 1 & 0 \\ 0 & -4/13 & 1 \end{bmatrix} \begin{bmatrix} 2 & 3 & 1 \\ 0 & -13/2 & 1/2 \\ 0 & 0 & 54/13 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} \Leftrightarrow$$

$$\begin{bmatrix} 2 & 3 & 1 \\ 0 & -13/2 & 1/2 \\ 0 & 0 & 54/13 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -2 \\ 6 \\ 24/13 \end{bmatrix}$$

Direct LU Solver Accuracy

Accuracy: $\|x - x^*\| / \|x^*\| = \phi(n) \kappa(A) \epsilon_0$

What is condition number $\kappa(A)$?

Sensitivity of Solution to perturbation: $1 \leq \kappa(A) < \infty$

$\epsilon_0 \Rightarrow 2^{-24}$ single precision, 2^{-53} double precision

Accuracy of LU solver is not enough

Iterative Refinement with LU Solver

Accuracy: $\|x - x^*\| / \|x^*\| = \phi(n) \kappa(A) \epsilon_0$

What is condition number $\kappa(A)$?

Sensitivity of the matrix A to perturbation: $1 \leq \kappa(A) < \infty$

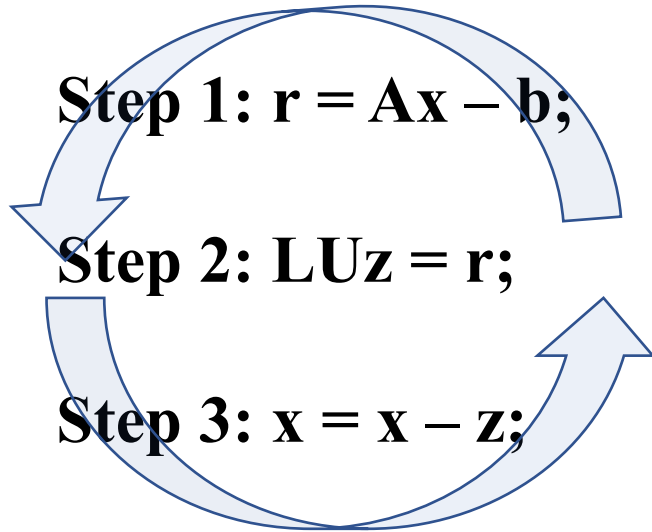
$\epsilon_0 \Rightarrow 2^{-24}$ single precision, 2^{-53} double precision

Accuracy of LU solver is not enough

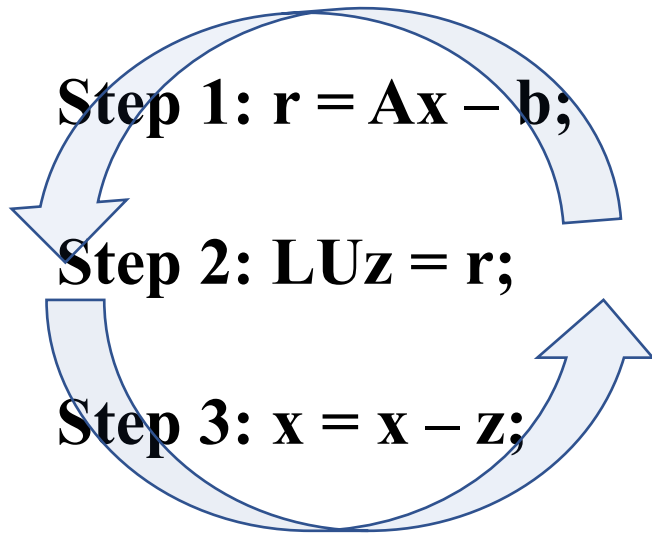
\Rightarrow Use it as an approximator inside Iterative Refinement

Iterative Refinement(IR): $q = \|x - x^*\| / \|x^*\|$, $q < 1$ for IR Success

Iterative Refinement with LU Solver



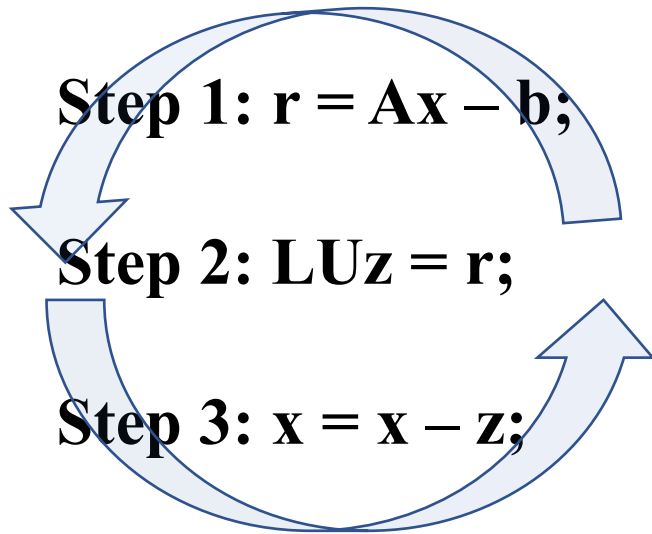
Iterative Refinement with LU Solver



Step 1:

Seeks r : $r = A(\Delta x)$; Matrix times previous solution error sought

Iterative Refinement with LU Solver



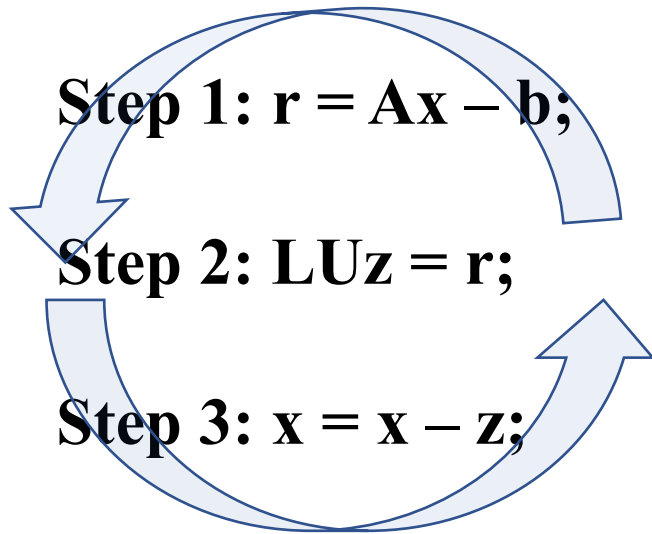
Step 1:

Seeks \mathbf{r} : $\mathbf{r} = \mathbf{A}(\Delta\mathbf{x})$; Matrix times previous solution error sought

Step 2:

Seeks \mathbf{z} : $\mathbf{z} = \mathbf{A}^{-1}\mathbf{A}(\Delta\mathbf{x})$; previous solution error approximation sought

Iterative Refinement with LU Solver



Step 1:

Seeks r : $r = A(\Delta x)$; Matrix times previous solution error sought

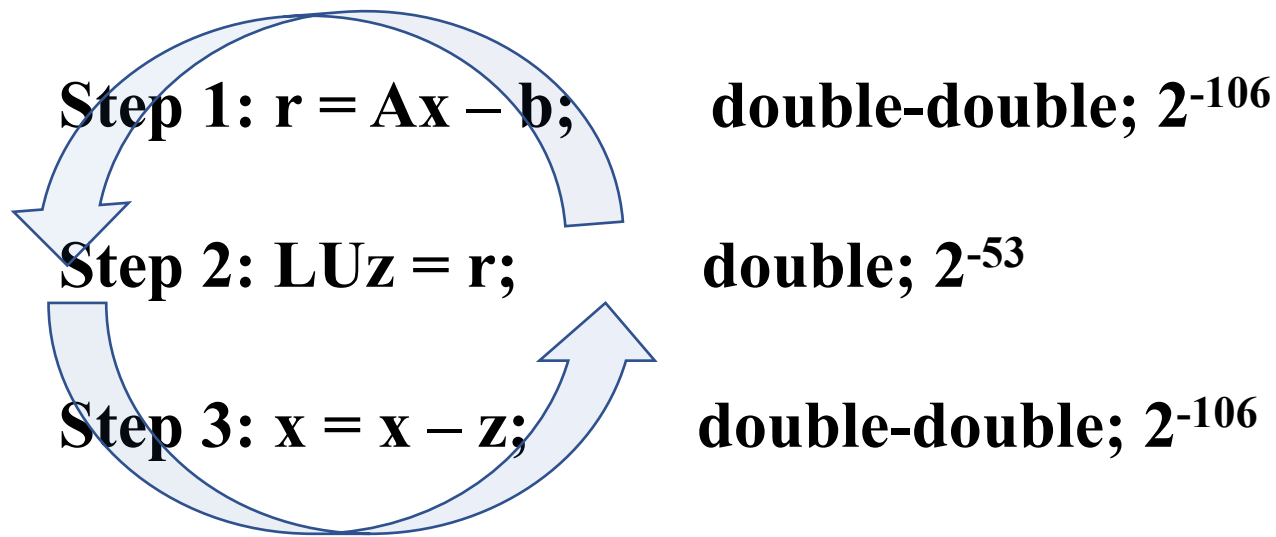
Step 2:

Seeks z : $z = A^{-1}A(\Delta x)$; previous solution error approximation sought

Step 3:

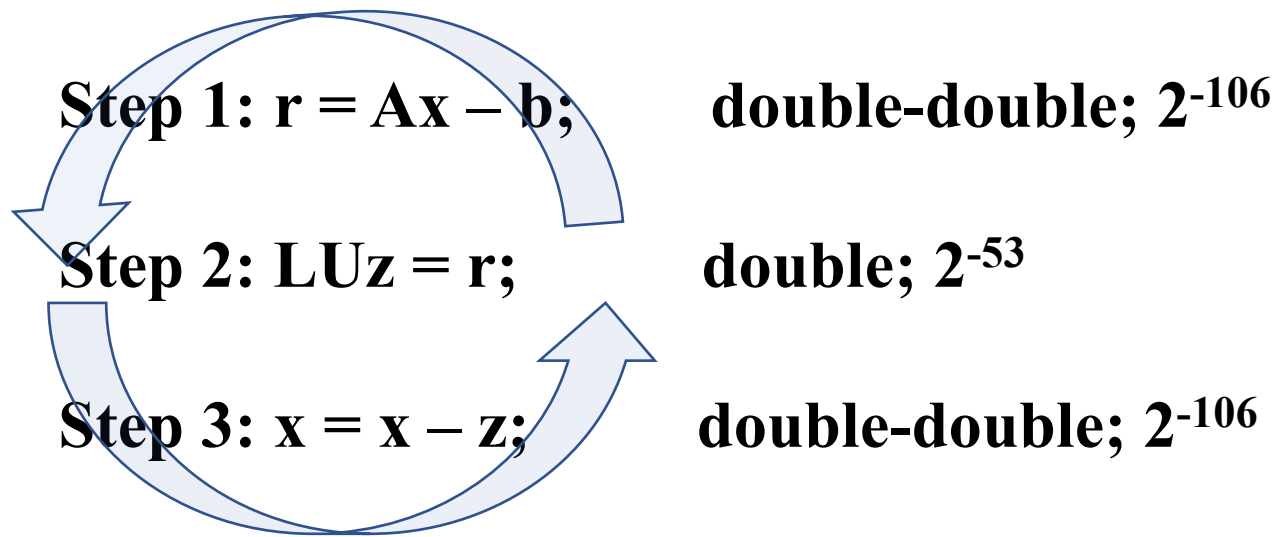
Update x : $x = x - (\Delta x)$; previous solution error approximation subtracted

Conventional Iterative Refinement



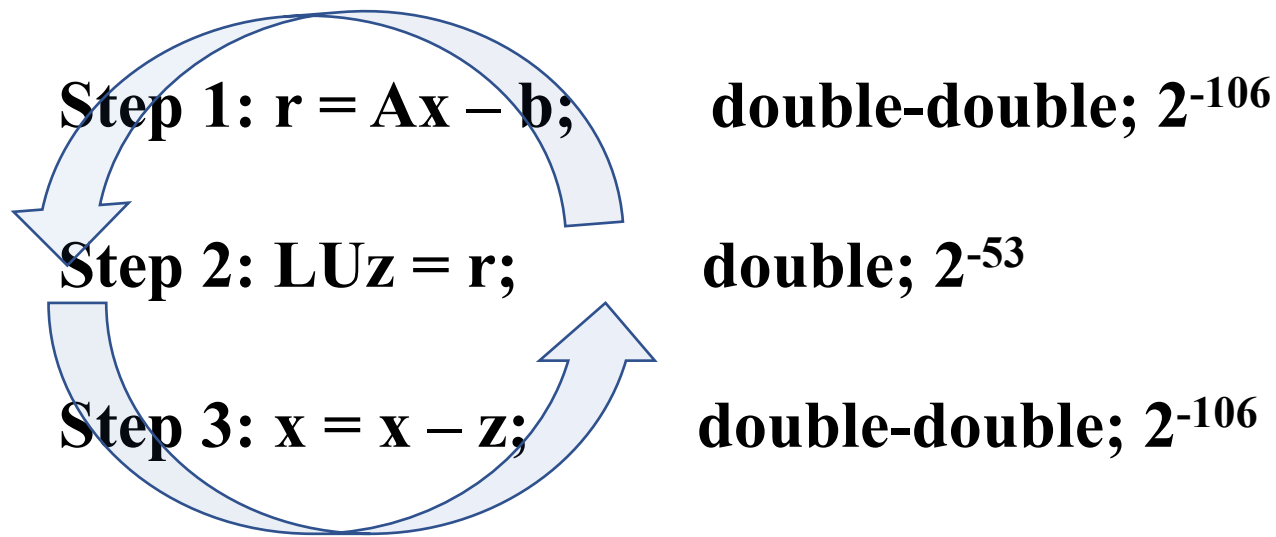
Accuracy: $\|x - x^*\| / \|x^*\| < 2^{-53}$

Conventional Iterative Refinement



Accuracy: $\|x - x^*\| / \|x^*\| < 2^{-53} \sim 10^{-16}$

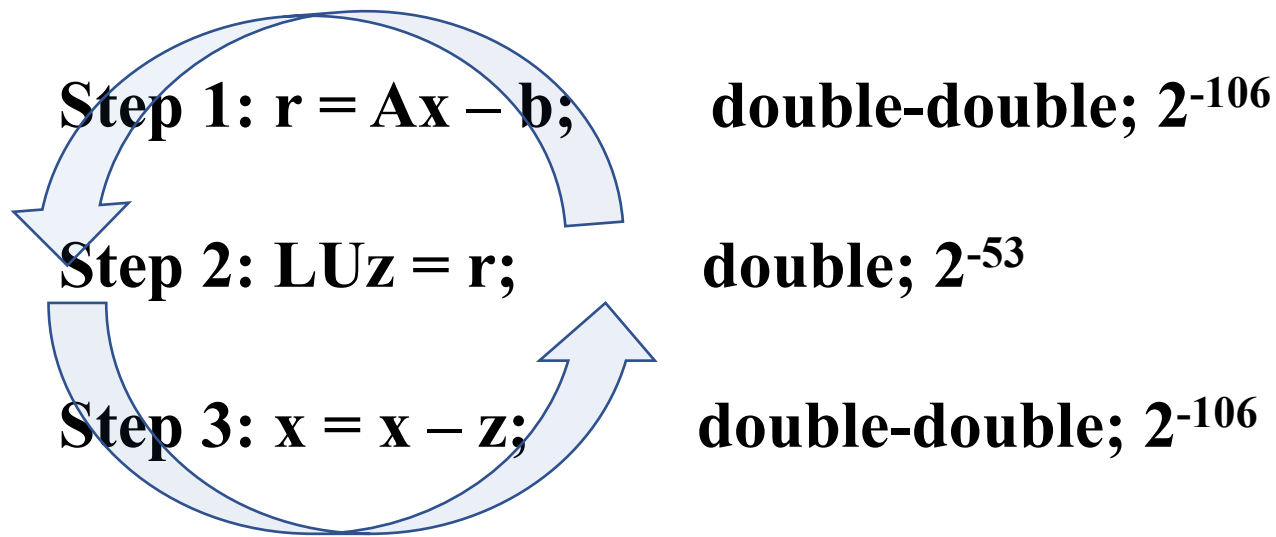
Conventional Iterative Refinement



Accuracy: $\|x - x^*\| / \|x^*\| < 2^{-53}$

Can we apply reduced precision arithmetic without losing accuracy? -> Transprecision Technique

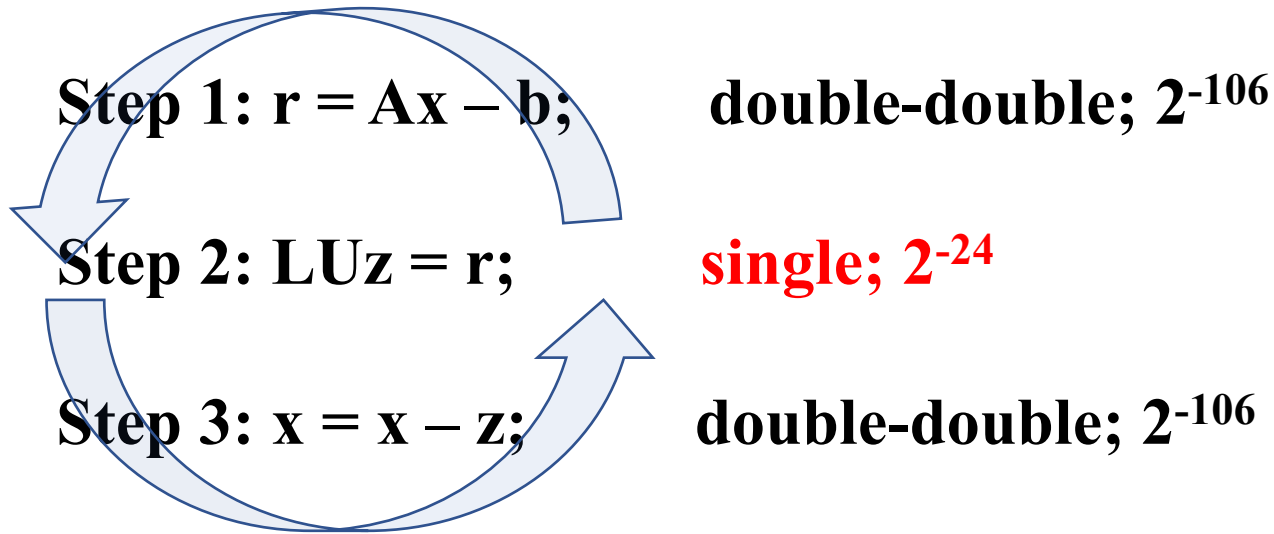
Conventional Iterative Refinement



Accuracy: $\|\mathbf{x} - \mathbf{x}^*\| / \|\mathbf{x}^*\| < 2^{-53}$

TT1: $\|\mathbf{z} - \mathbf{z}^*\| / \|\mathbf{z}^*\| < 1$ for **LU solver** in Step 2.

Mixed Precision Iterative Refinement



Accuracy: $\|\mathbf{x} - \mathbf{x}^*\| / \|\mathbf{x}^*\| < 2^{-53}$

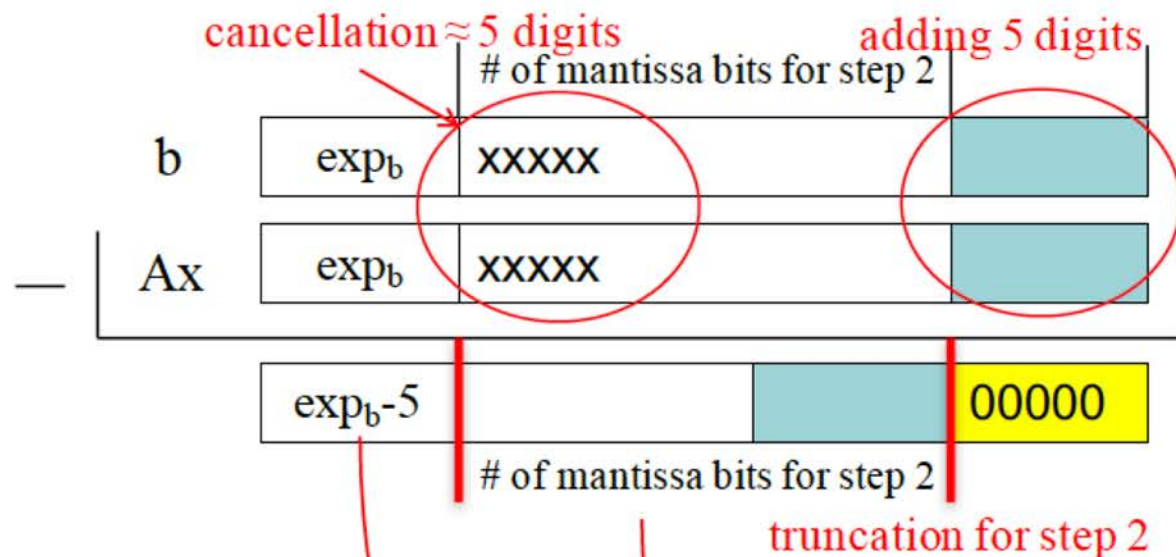
TT1: $\|\mathbf{z} - \mathbf{z}^*\| / \|\mathbf{z}^*\| < 1$ for LU solver in Step 2.

=> This is generally referred to Mixed Precision IR.

Transprecision Iterative Refinement

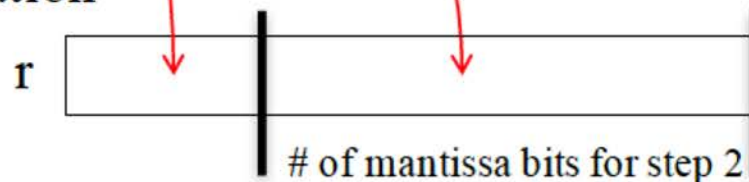
seeking residual

step 1 $r^{(i)} = Ax^{(i)} - b$

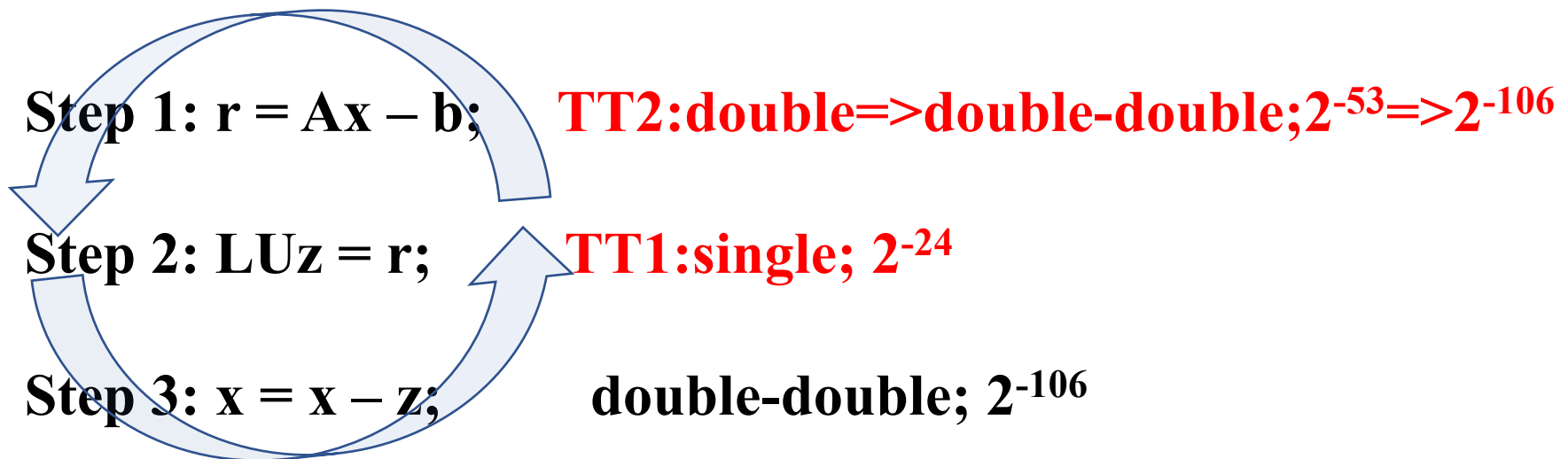
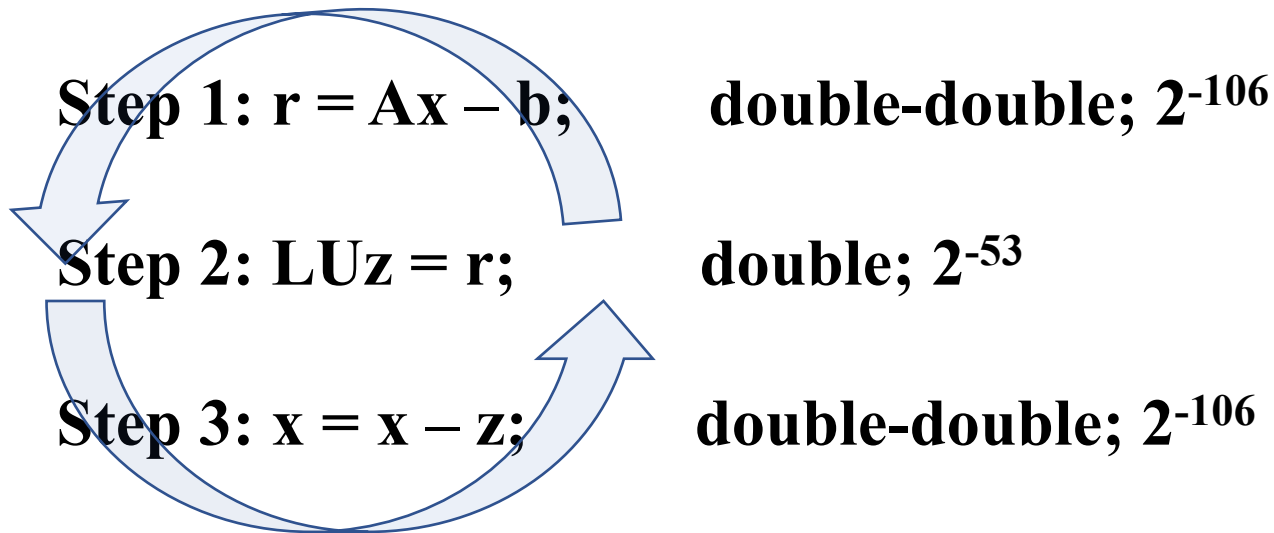


seeking the error at the previous iteration

step 2 $Az^{(i)} = r^{(i)}$



Transprecision Iterative Refinement



Transprecision Iterative Refinement

Beyond this lecture, other two transprecision techniques are employed.

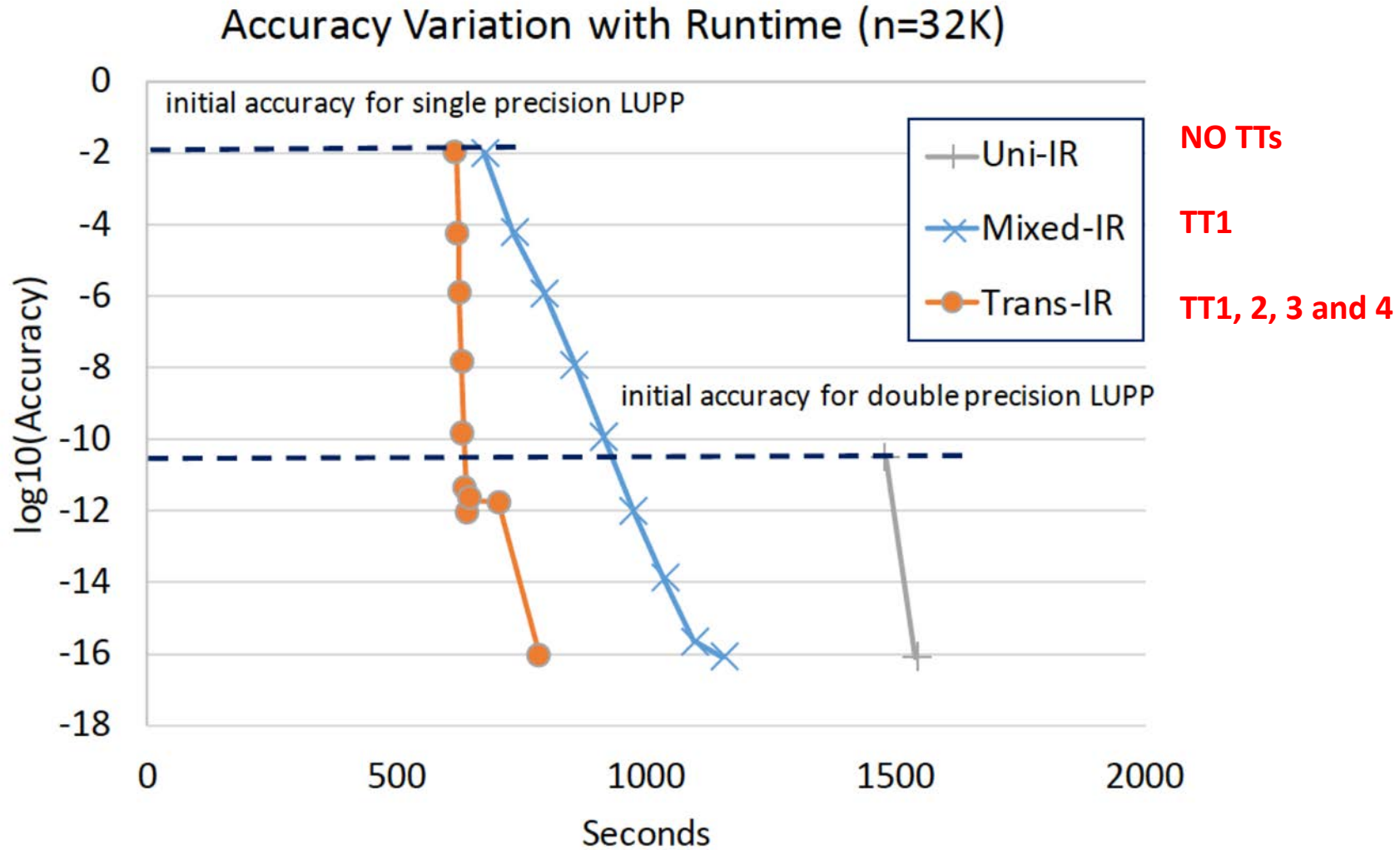
TT3: Recursive iterative refinement

TT4: Skipping accuracy checks

Reference for details of all TTs:

J. Lee, H. Vandierendonck, M. Arif, G. D. Peterson and D. S. Nikolopoulos, "Energy-Efficient Iterative Refinement using Dynamic Precision," in *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*. doi: 10.1109/JETCAS.2018.2850665

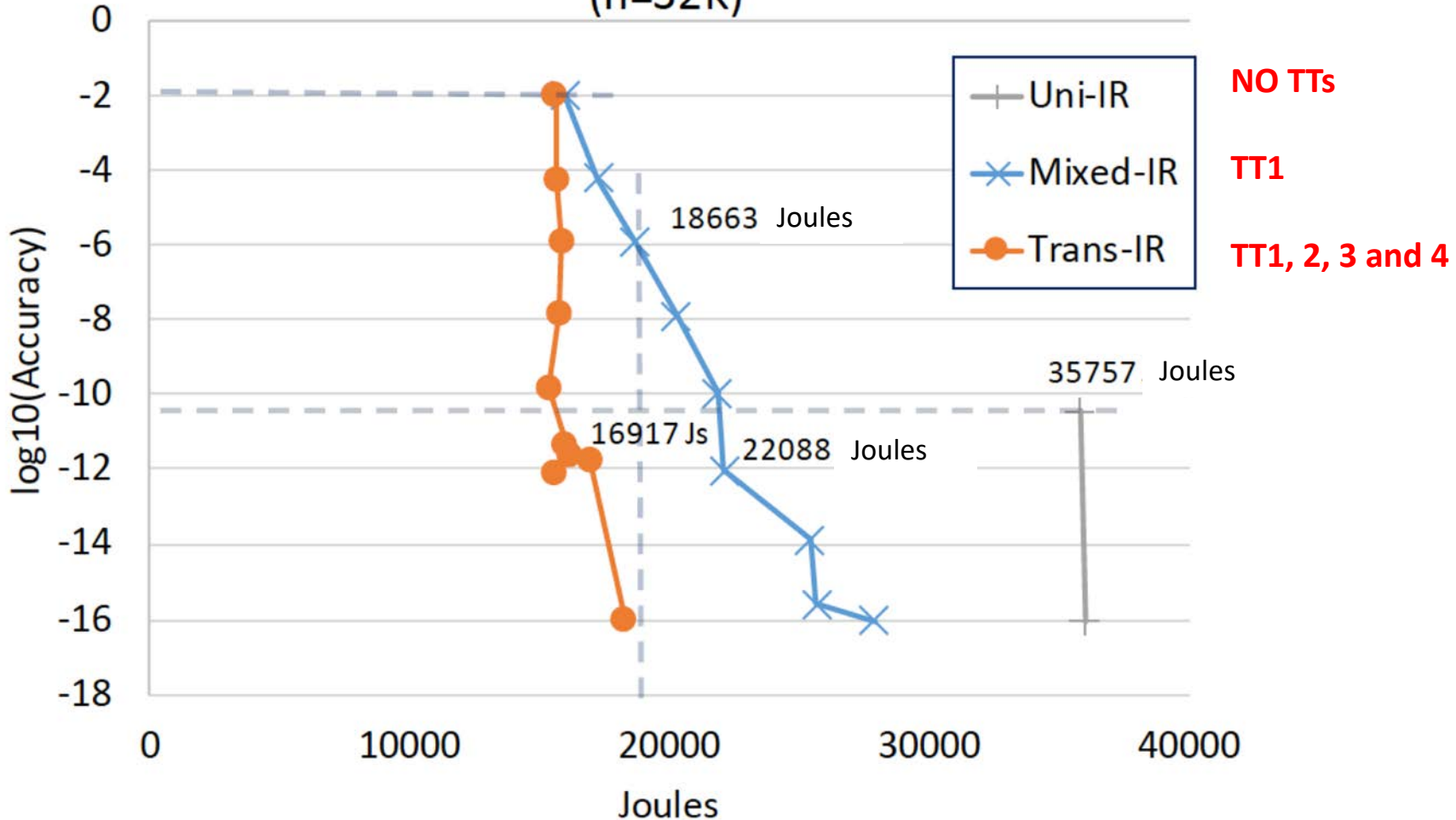
Impact of Transprecision Iterative Refinement



Impact of Transprecision Iterative Refinement

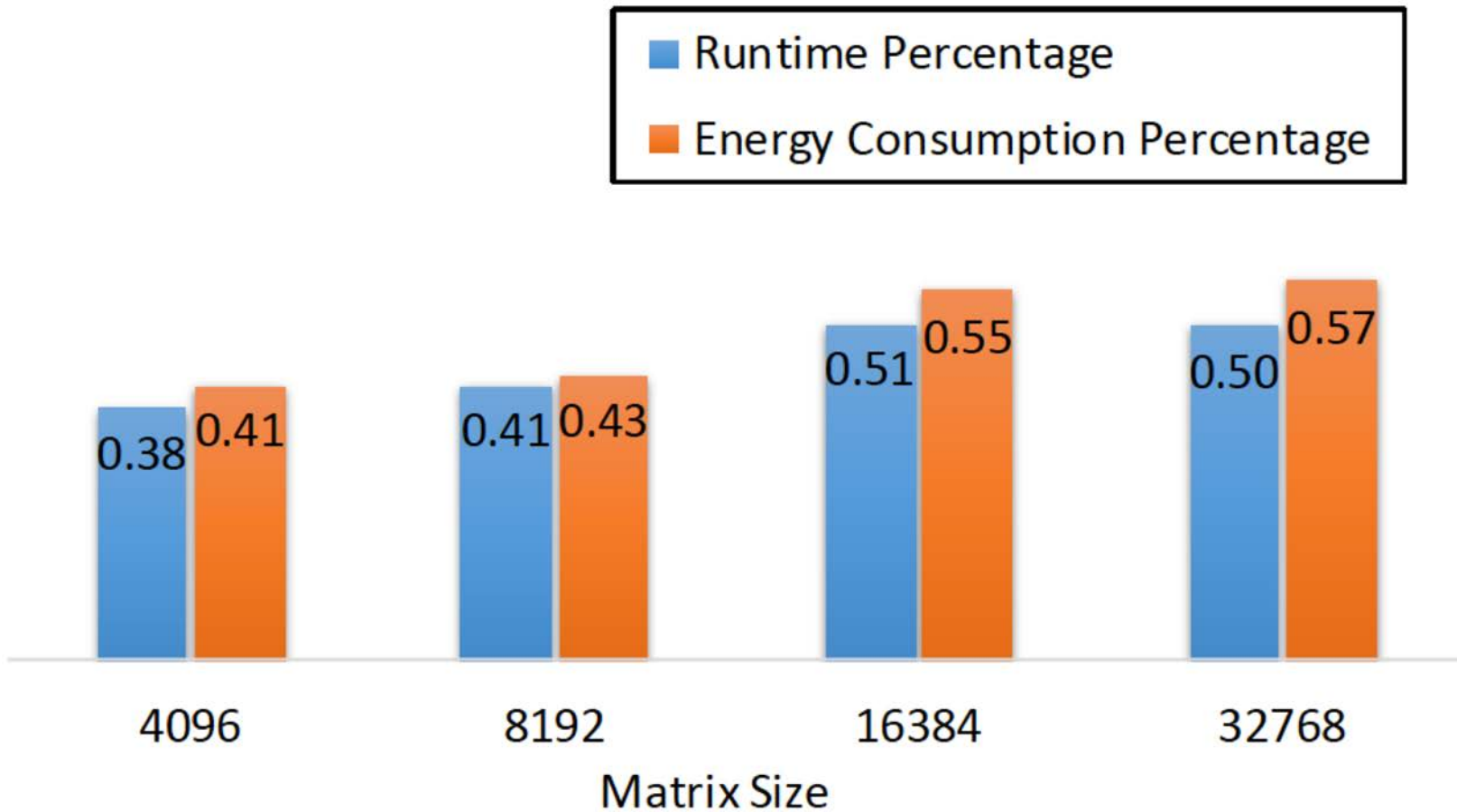
Accuracy Variation with Energy Consumption

(n=32K)



Impact of Transprecision Iterative Refinement

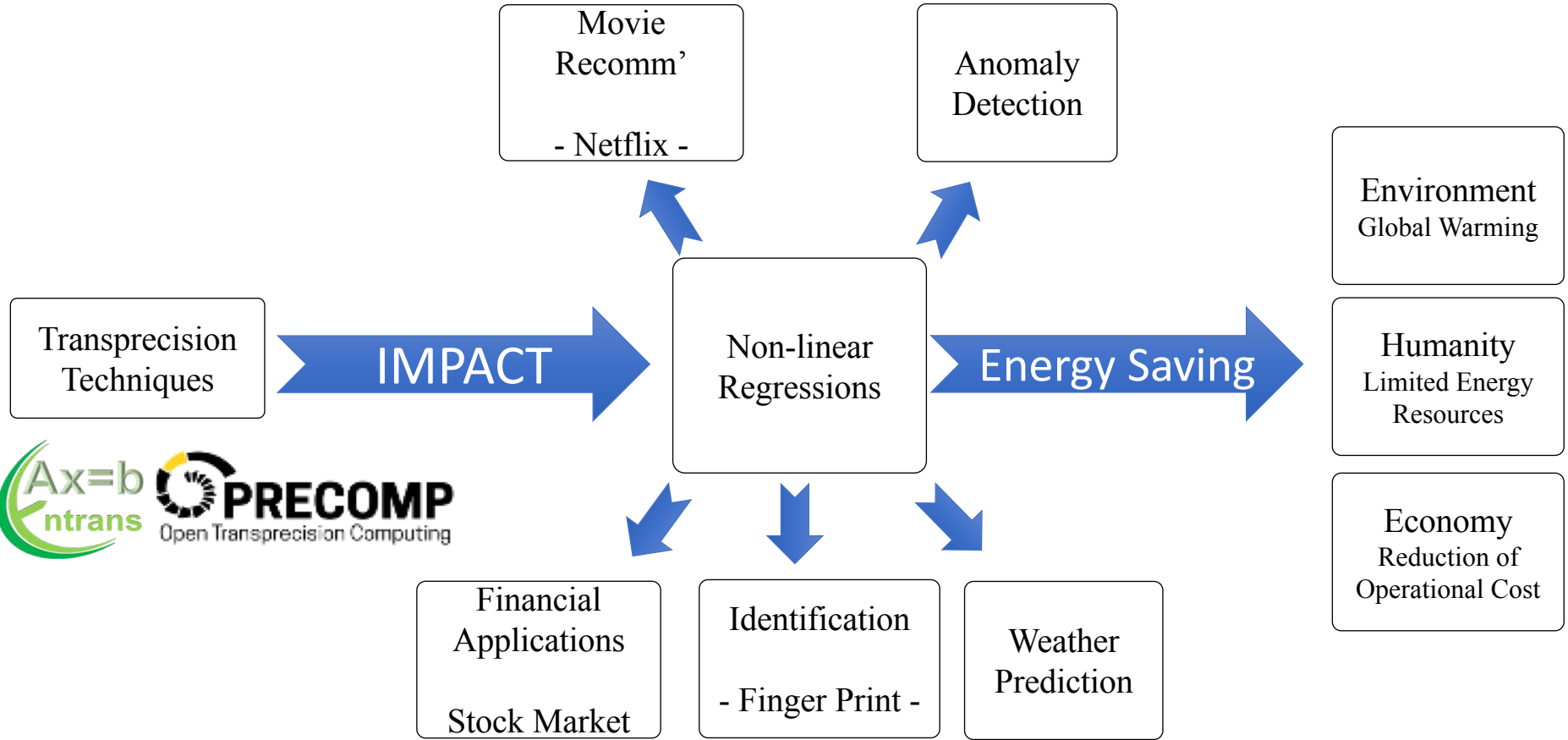
Percentages (baseline: Mixed-IR) **TT1 vs TT1, 2, 3 and 4**



Part B. Non-linear Regressions



Impact of Transprecision on Non-linear Regression Applications



Linear Regression and Non-linear Regression

Linear Regression (e.g., Recursive Least Squares) :

Seek $\mathbf{w} = (w_0, w_1, w_2, w_3)$ to estimate y given an input $\mathbf{x} = (1, x_1, x_2, x_3)$ with

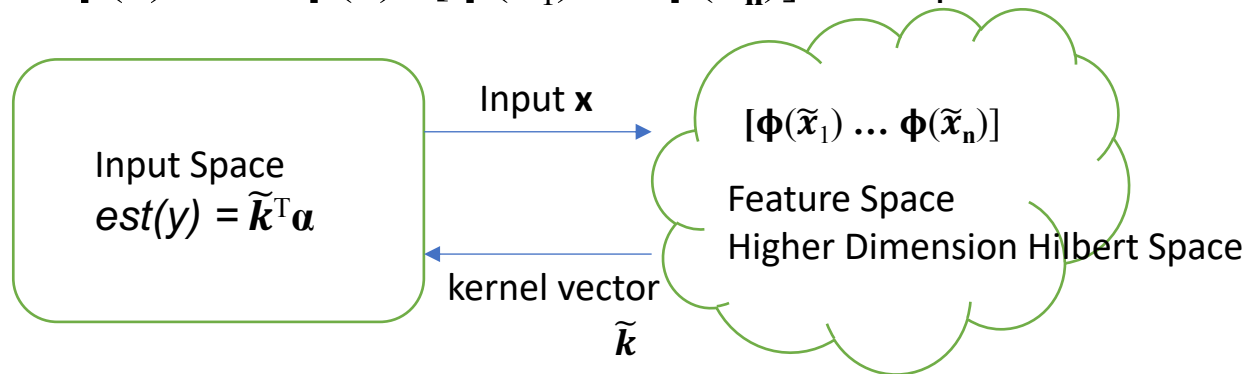
$$\text{est}(y) = \mathbf{x}^T \mathbf{w} = w_0 \cdot 1 + w_1 x_1 + w_2 x_2 + w_3 x_3.$$

Non-linear Regression using **Kernel Method**

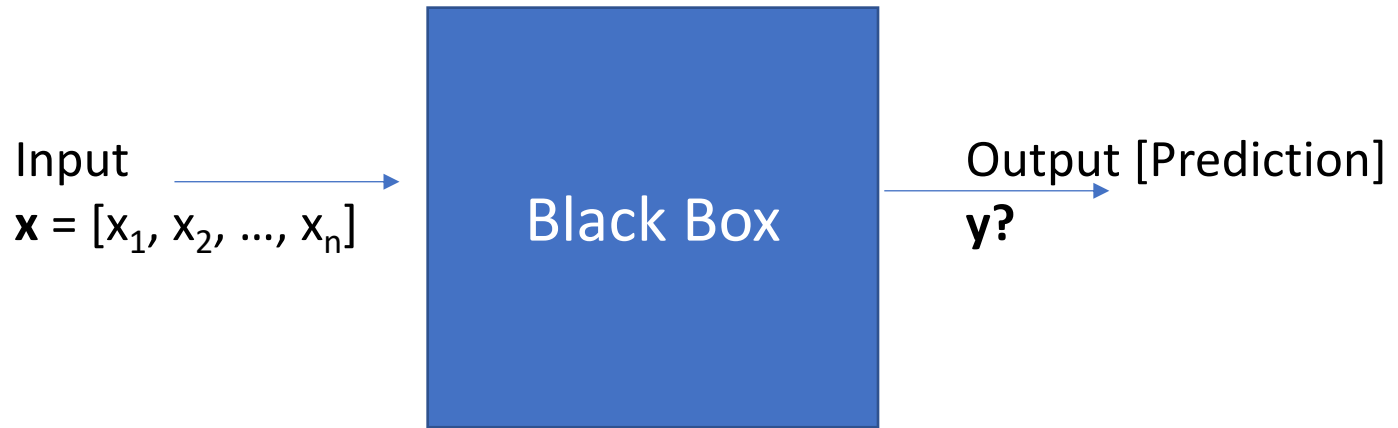
(e.g., Kernel Recursive Least Squares):

Perform regression in **higher dimension Hilbert space \mathcal{H}**

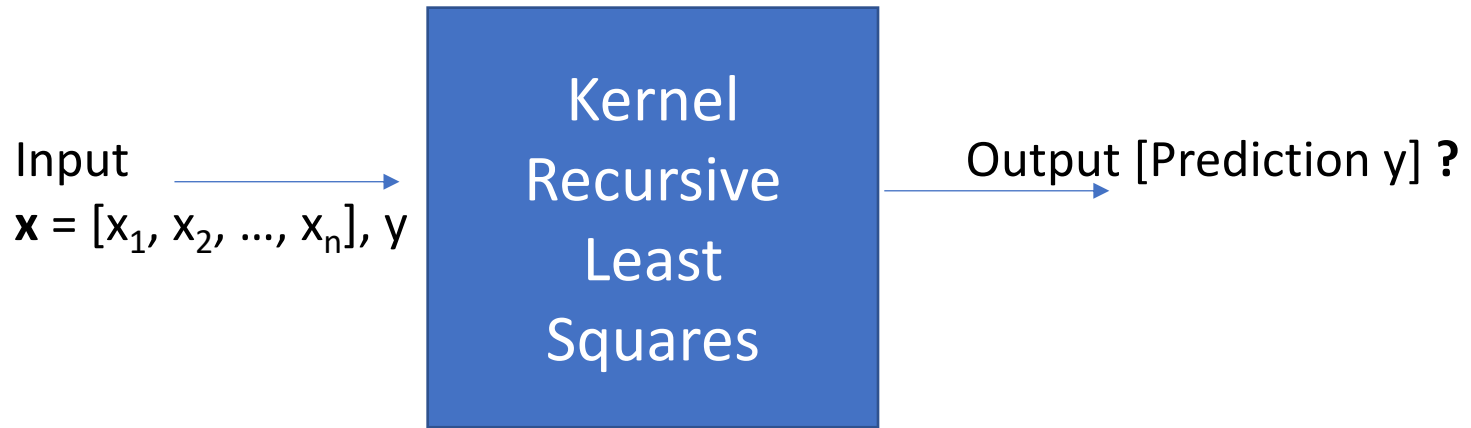
$$\text{est}(y) = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{w} = \boldsymbol{\phi}(\mathbf{x})^T [\boldsymbol{\phi}(\tilde{\mathbf{x}}_1) \dots \boldsymbol{\phi}(\tilde{\mathbf{x}}_n)] \boldsymbol{\alpha} = \tilde{\mathbf{k}}^T \boldsymbol{\alpha}$$



Non-linear Regression Machine Learning

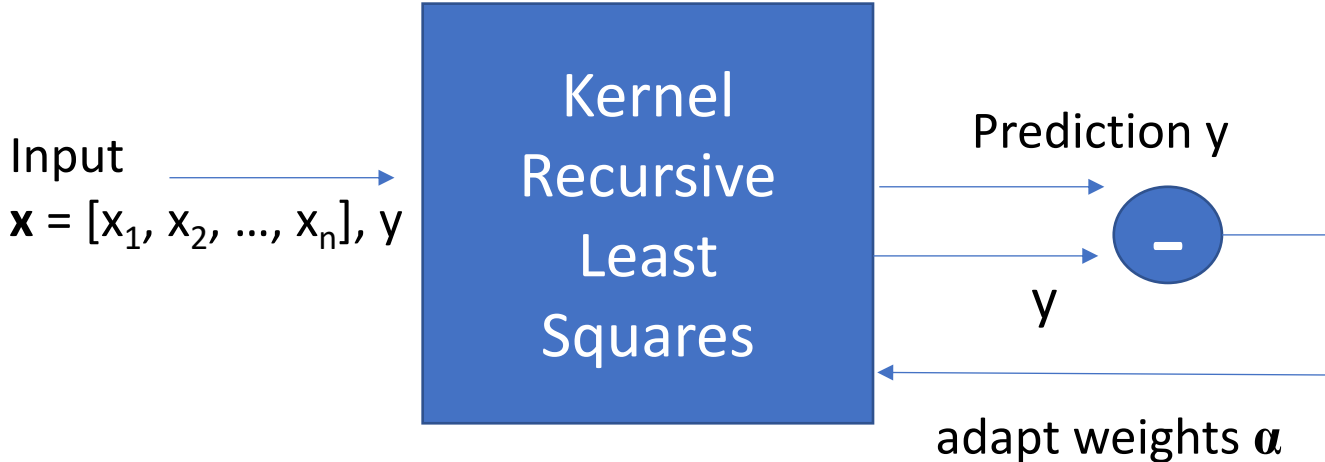


Non-linear Regression Machine Learning



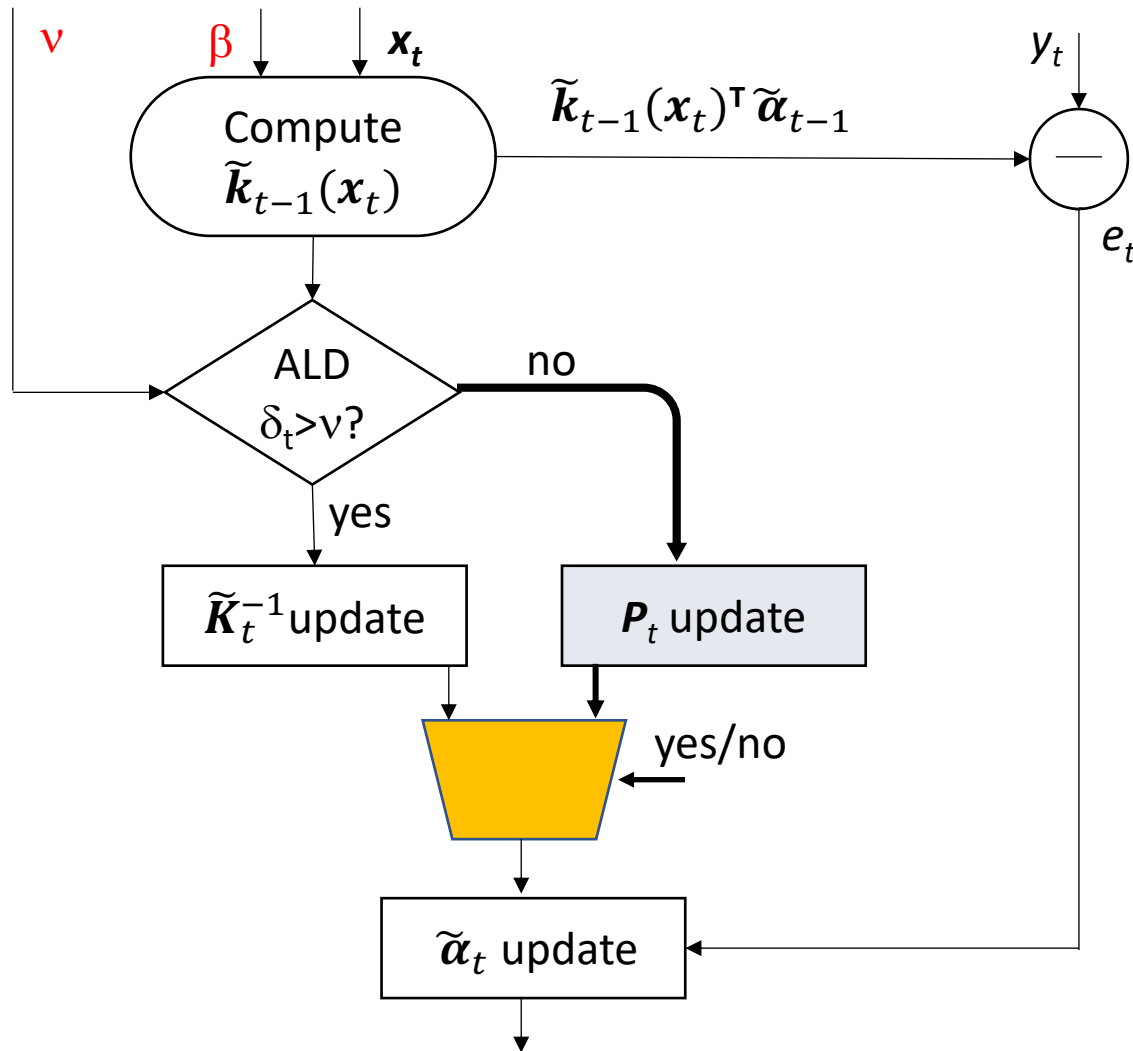
Realtime Training

Non-linear Regression Machine Learning



Realtime Training
(i.e., Kernel Adaptive Filter)

KRLS Algorithm Computing Components



Kernel Recursive Least Squares (KRLS)

Applications: Non-linear regressions

- Simple RLS mechanism
- No local minima issue
- Fast Convergence / Good Prediction Accuracy
- Online Adaptive Learning (Learning weights One by One – Fit for large scale Machine Learning)

Model Selection for KRLS

KRLS model depending on two hyperparameters :

Approx Linear Dependency(ALD) threshold ν , Kernel Width β

ALD threshold ν : Generalization and Approximation level

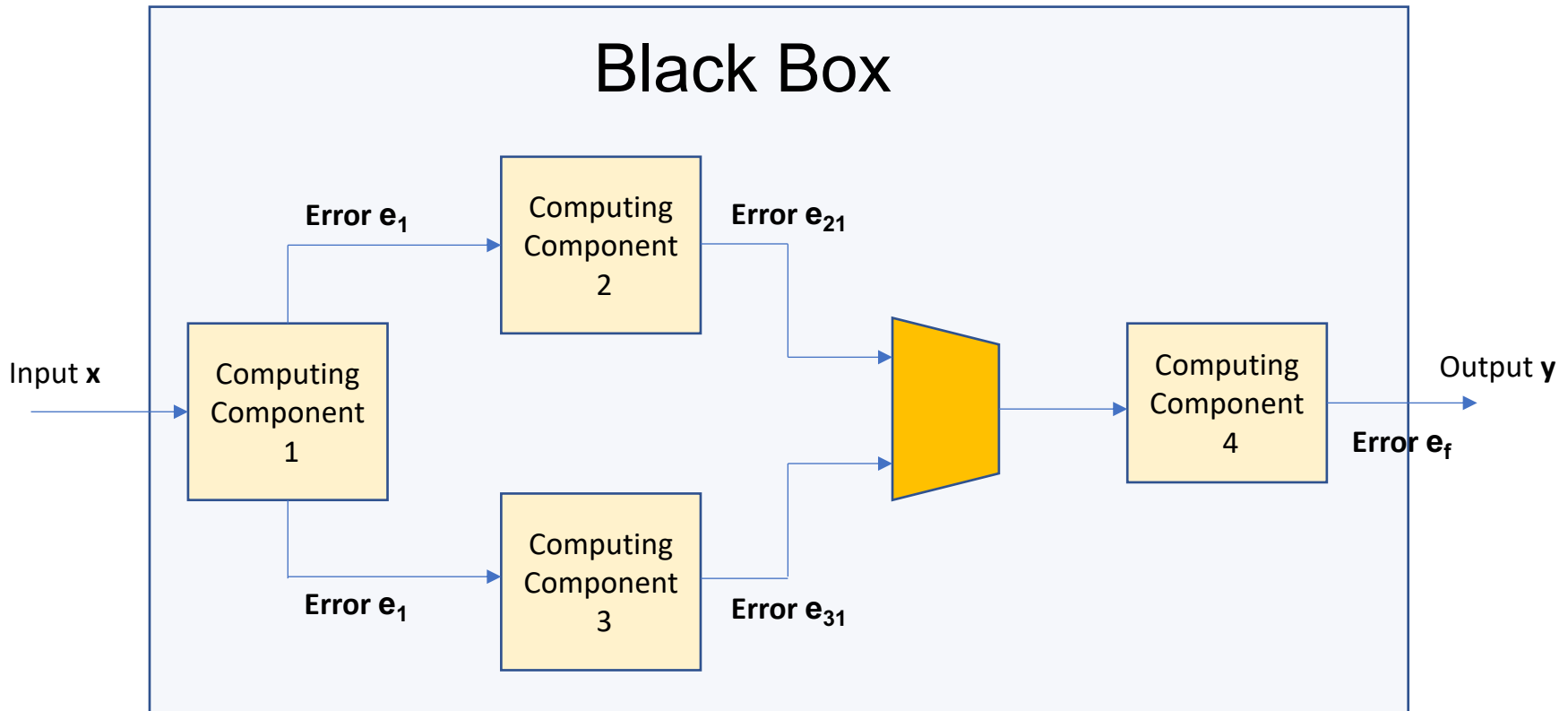
Kernel Width β : Inner product distribution

$$\kappa(\mathbf{x}_t, \mathbf{x}_i) = \langle \phi(\mathbf{x}_t), \phi(\mathbf{x}_i) \rangle = \exp^{-\frac{(\mathbf{x}_t - \mathbf{x}_i)^T (\mathbf{x}_t - \mathbf{x}_i)}{2\beta^2}}$$

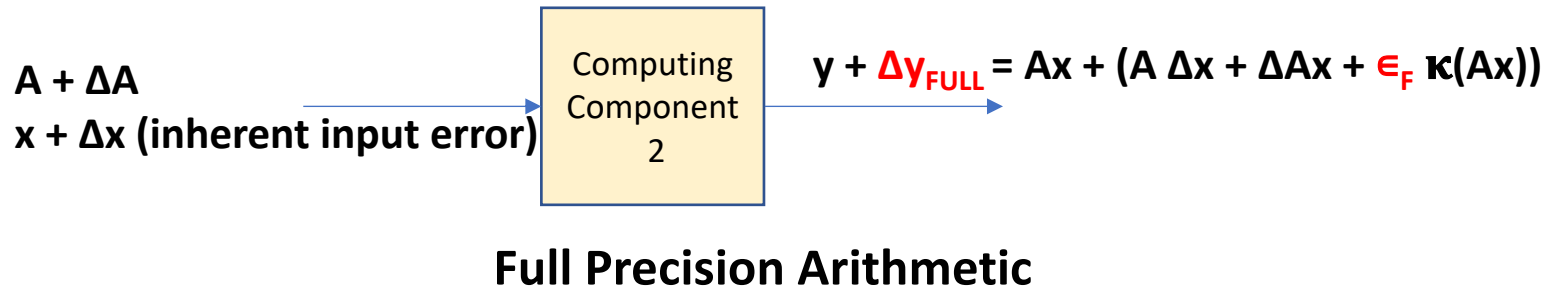
Optimal ν, β depends on Inputs

In many cases, Cross Validation decides Optimal ν, β

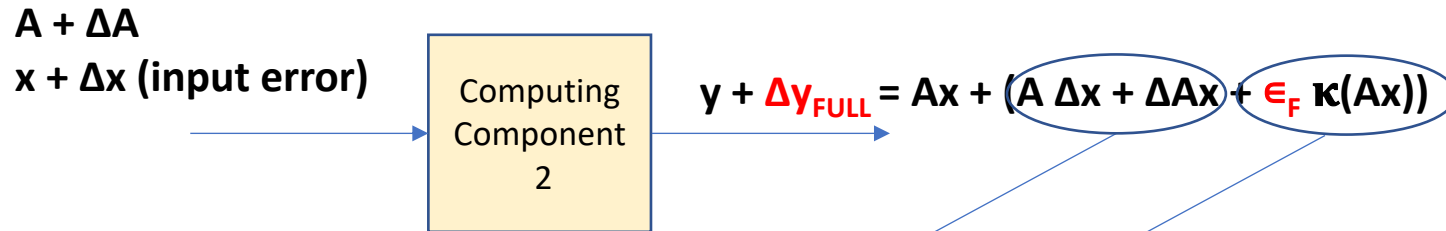
Error propagation for each computing component



Error propagation for each computing component



Error propagation for each computing component



Full Precision Arithmetic

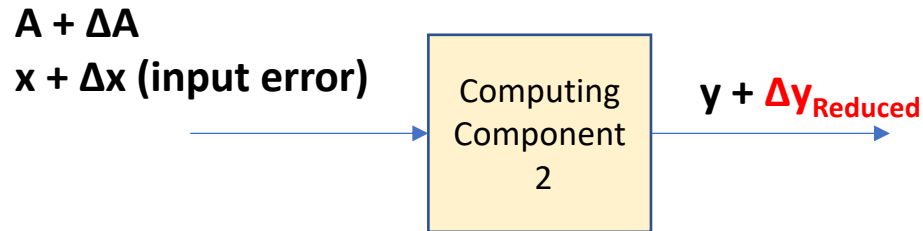
$||\Delta y_{\text{FULL}}||$ includes

err_{una} (Unavoidable) and err_{rnd} (Controllable: Precision arithmetic related)

Key Idea: ϵ can be balanced to err_{una} .

Employ lower precision arithmetic ϵ until $||\Delta y_{\text{FULL}}||$ not affected

Error propagation for each computing component



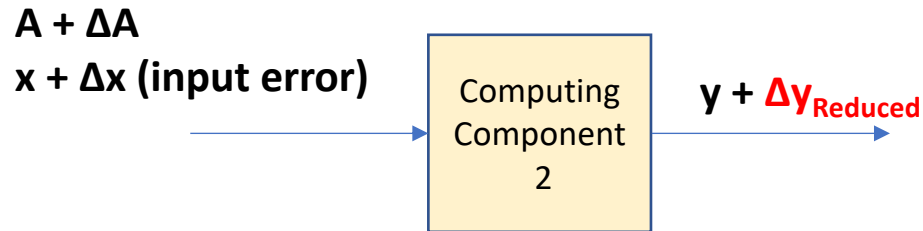
Reduced Precision Arithmetic

$$\text{err}_{\text{una}} \approx \text{err}_{\text{rnd}}$$

When size of Unavoidable Error equivalent of size of rounding off error

=> Reduced Arithmetic can be applied to the Computing Component

Error propagation for each computing component



Reduced Precision Arithmetic

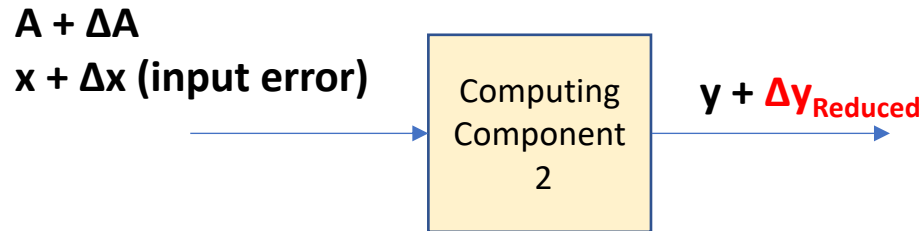
$$|| A\Delta x + \Delta Ax || \approx || \epsilon_R \kappa(Ax) ||$$

(unavoidable err)

(controllable rounding-off error)

Reduced Arithmetic can be applied to the Computing Component

Error propagation for each computing component



Reduced Precision Arithmetic

Case 1: Linear Solver

$$\text{err}_{\text{rnd}} \propto \kappa(A) \epsilon_R$$

Case 2: Matrix-vector

$$\text{err}_{\text{rnd}} \propto \kappa(A) \epsilon_R$$

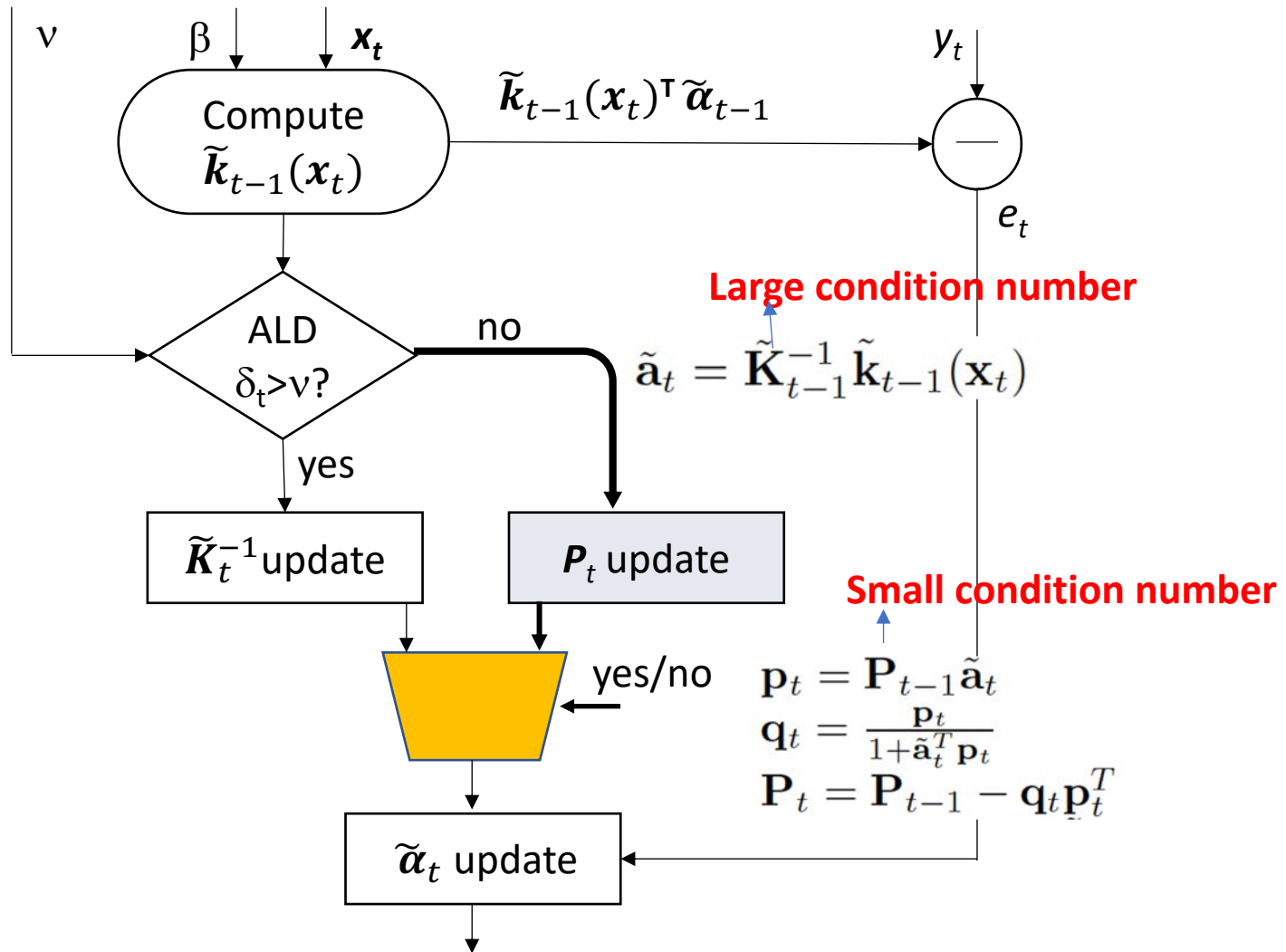
Case 3: Dot-product

$$\text{err}_{\text{rnd}} \propto (|\mathbf{x}_1|^T |\mathbf{x}_2|) / |\mathbf{x}_1^T \mathbf{x}_2| \epsilon_R$$

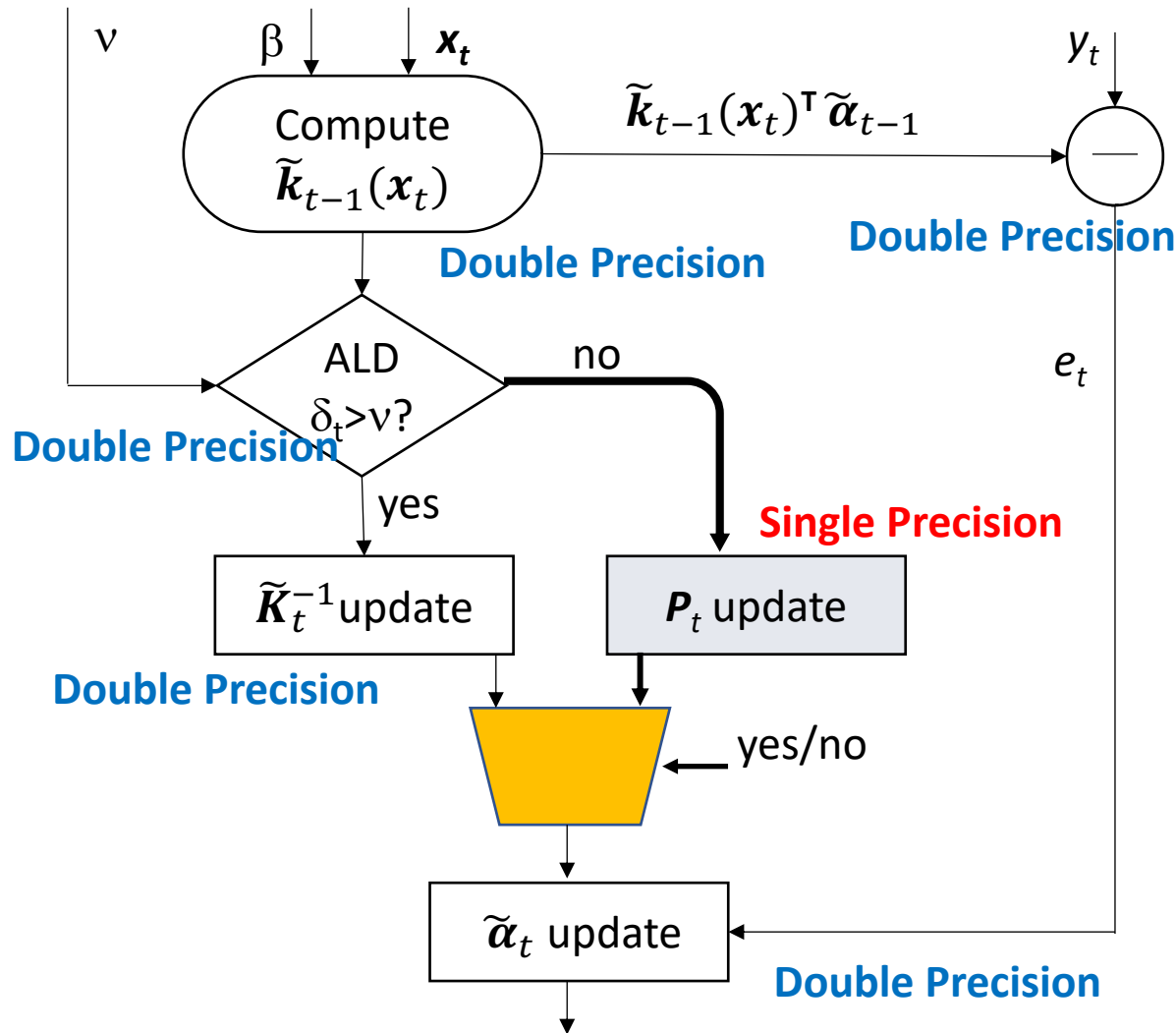
How can we use such properties?

Let us look at a Case Study with KRLS algorithm.

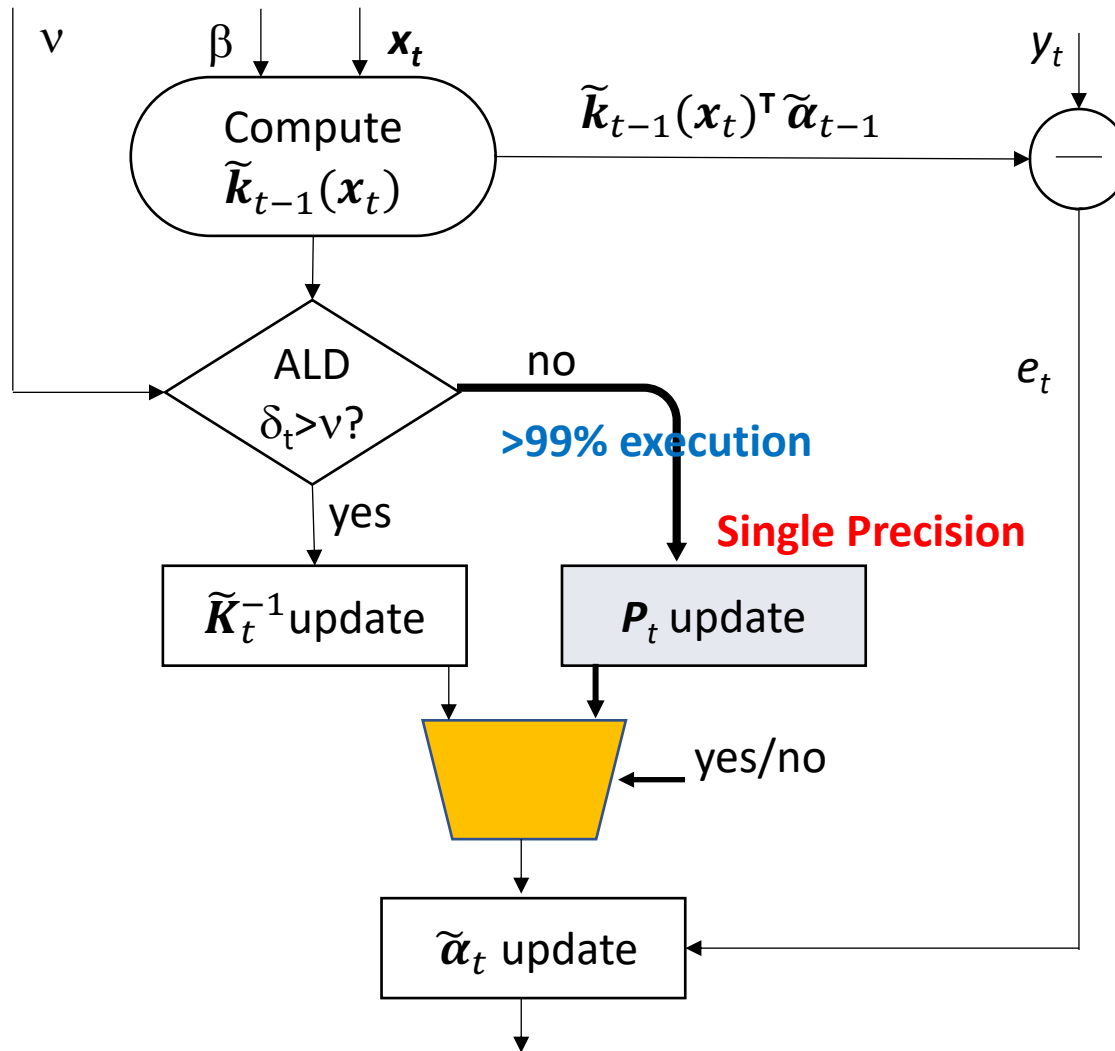
KRLS Algorithm Computing Components



Transprecision Computing for KRLS

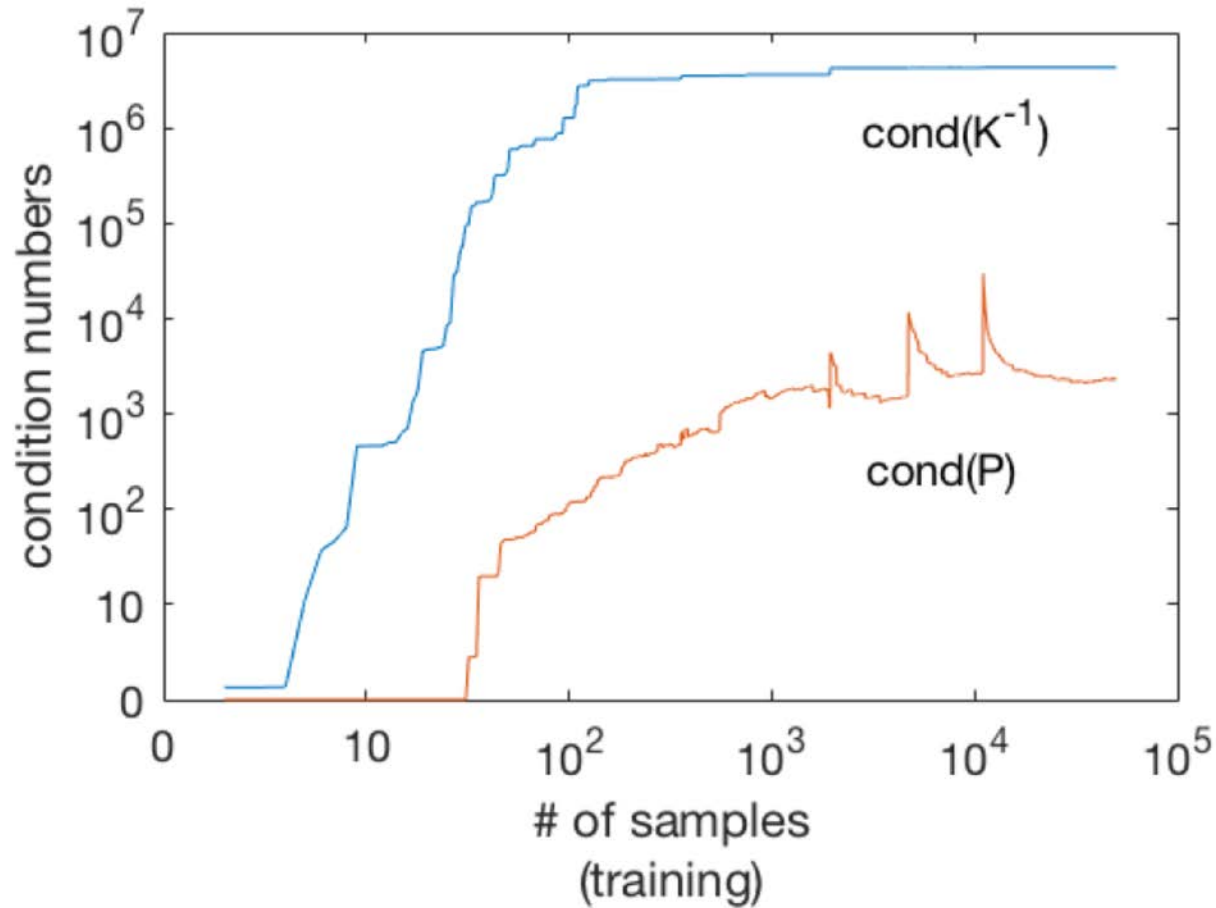


Transprecision Computing for KRLS



Condition Numbers of Matrices in KRLS

Condition number of K^{-1} and P with ALD v and kernel width β from Cross-validation



Impact of Transprecision KRLS

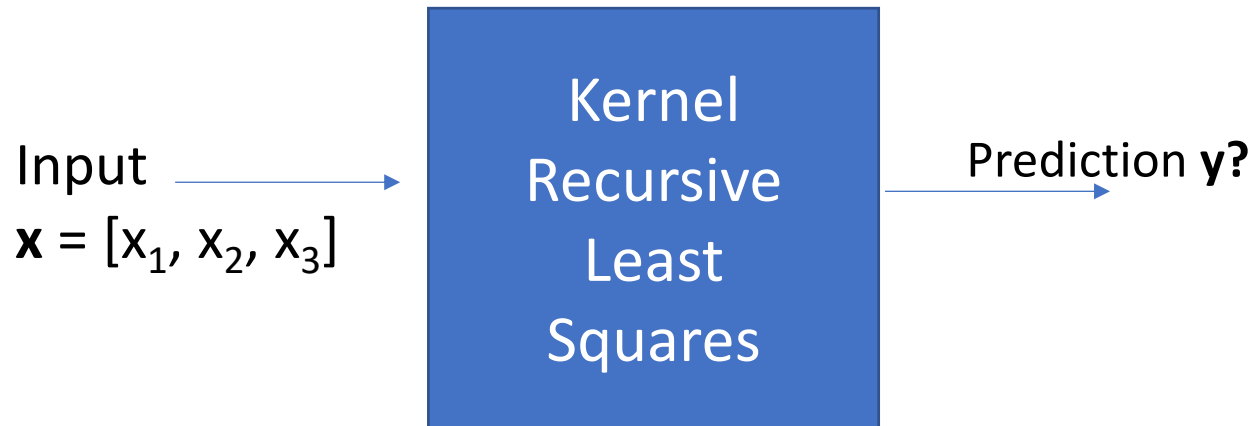
Non-Linear Regression Problem:

$$y = \sin(x_1)/x_1 + x_2/10.0 + \cos(x_3)$$

Impact of Transprecision KRLS

Non-Linear Regression Problem:

$$y = \sin(x_1)/x_1 + x_2/10.0 + \cos(x_3)$$



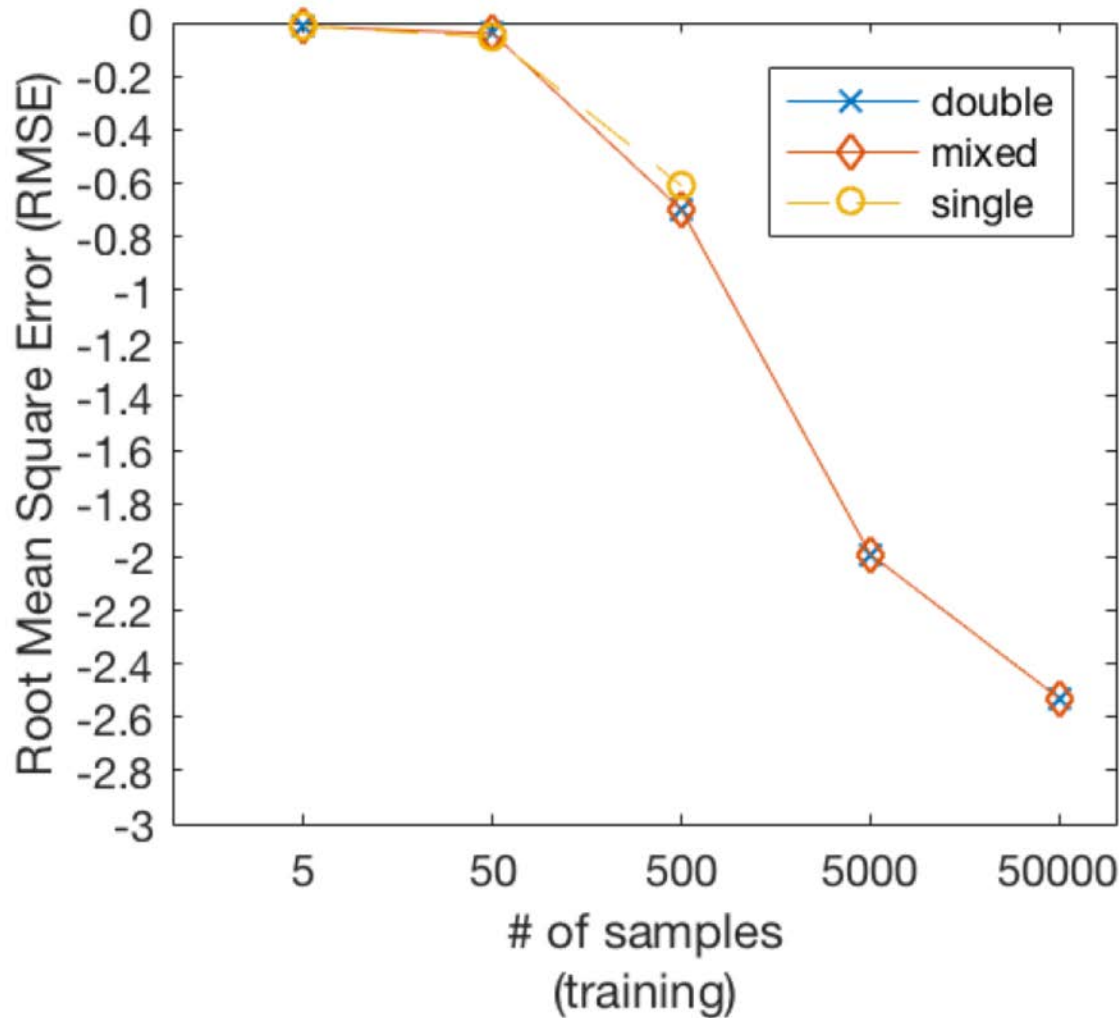
Impact of Transprecision KRLS

Non-Linear Regression Problem:

$$y = \sin(x_1)/x_1 + x_2/10.0 + \cos(x_3)$$

- $x_1, x_2, x_3 =$ Uniform Random $[-10, 10]$
- Gaussian Kernel Width $\beta = 2.5$
- ALD $\nu = 0.01$
- Intel(R) Xeon(R) CPU E5-2650 2GHz Single Core
- Energy Estimation: ALEA Energy Profiling Tool

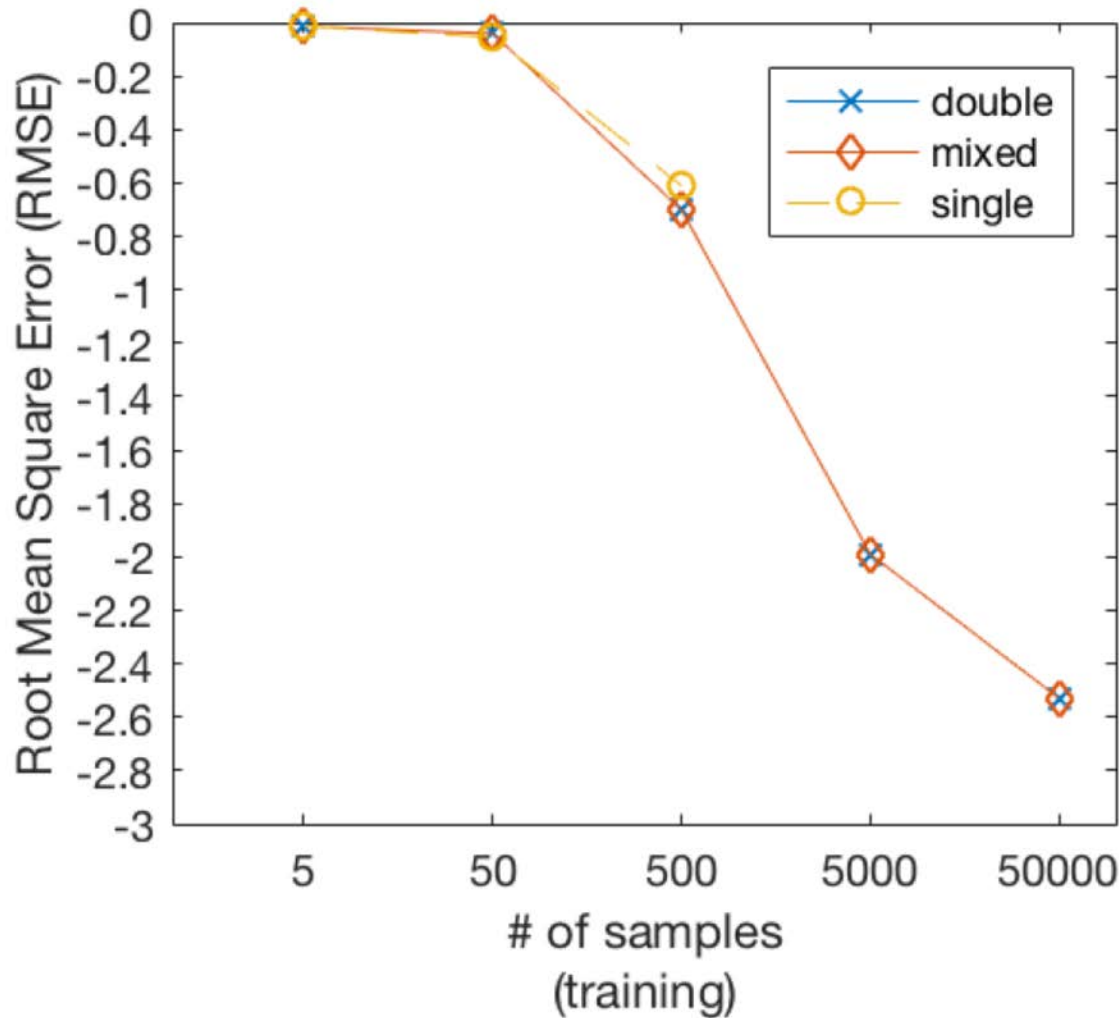
Prediction Accuracy



Almost 3 decimal digit accuracy when # training data = 50,000

RMSE is used for accuracy measure

Prediction Accuracy

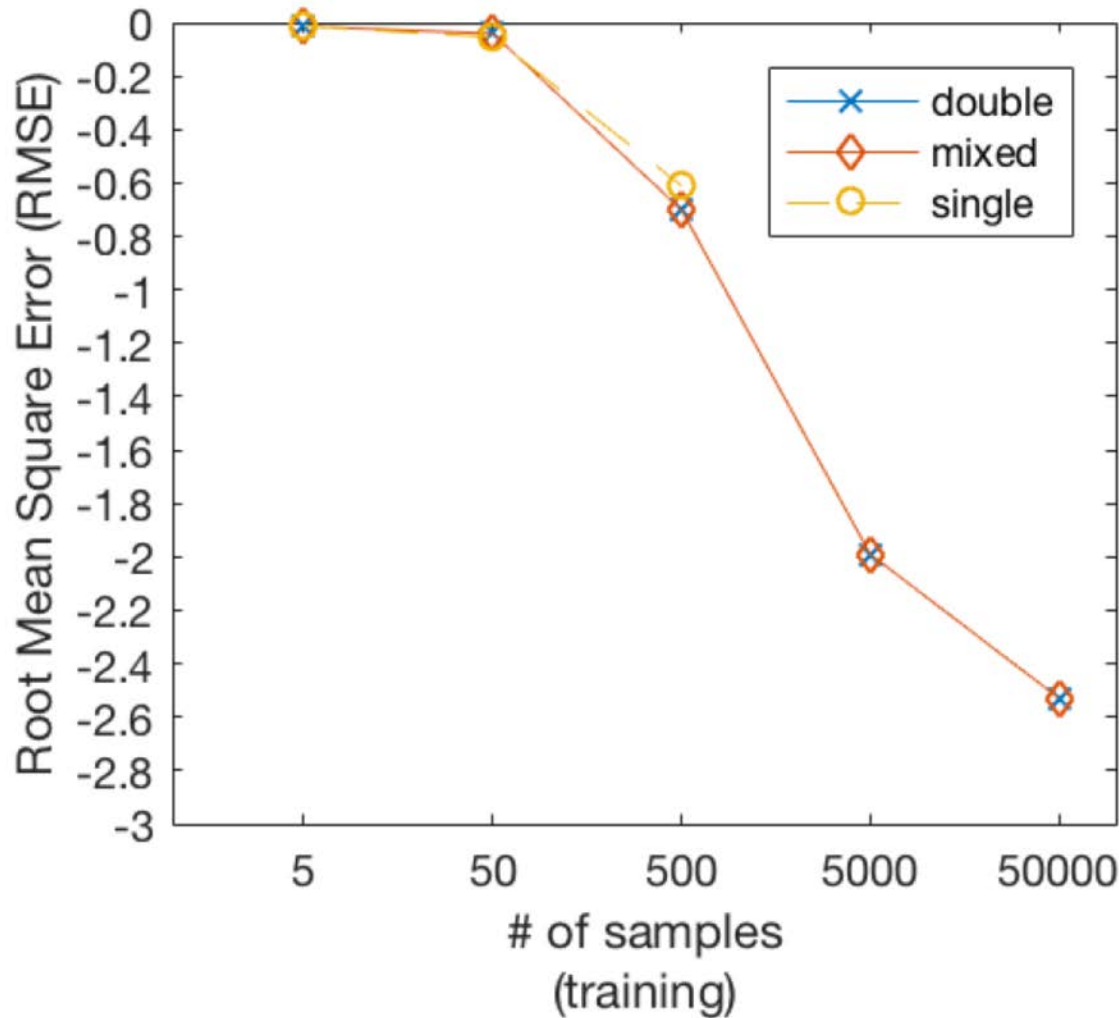


Almost 3 decimal digit accuracy when # training data = 50,000

Entire Single Precision **breaks** KRLS



Prediction Accuracy

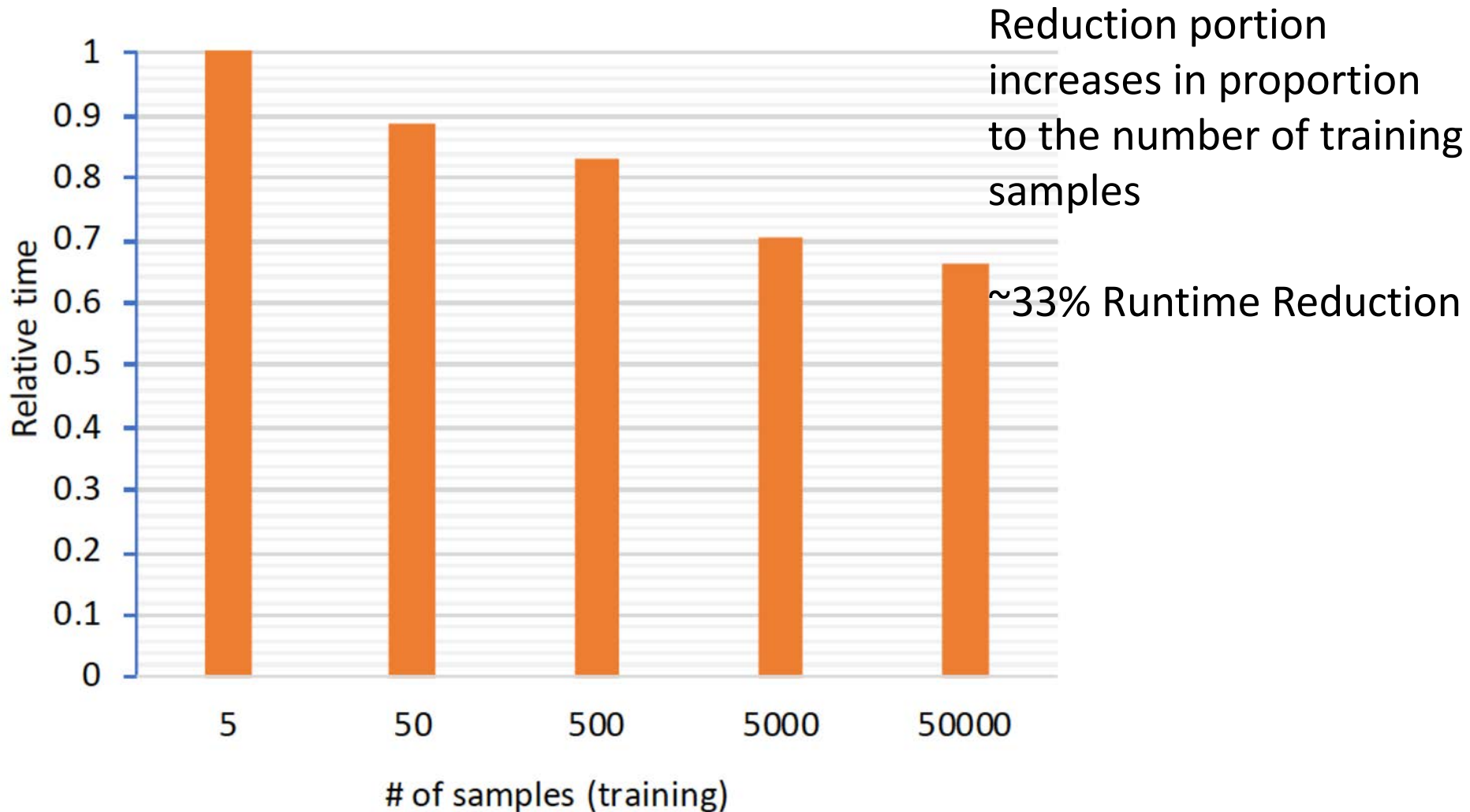


Almost 3 decimal digit accuracy when # training data = 50,000

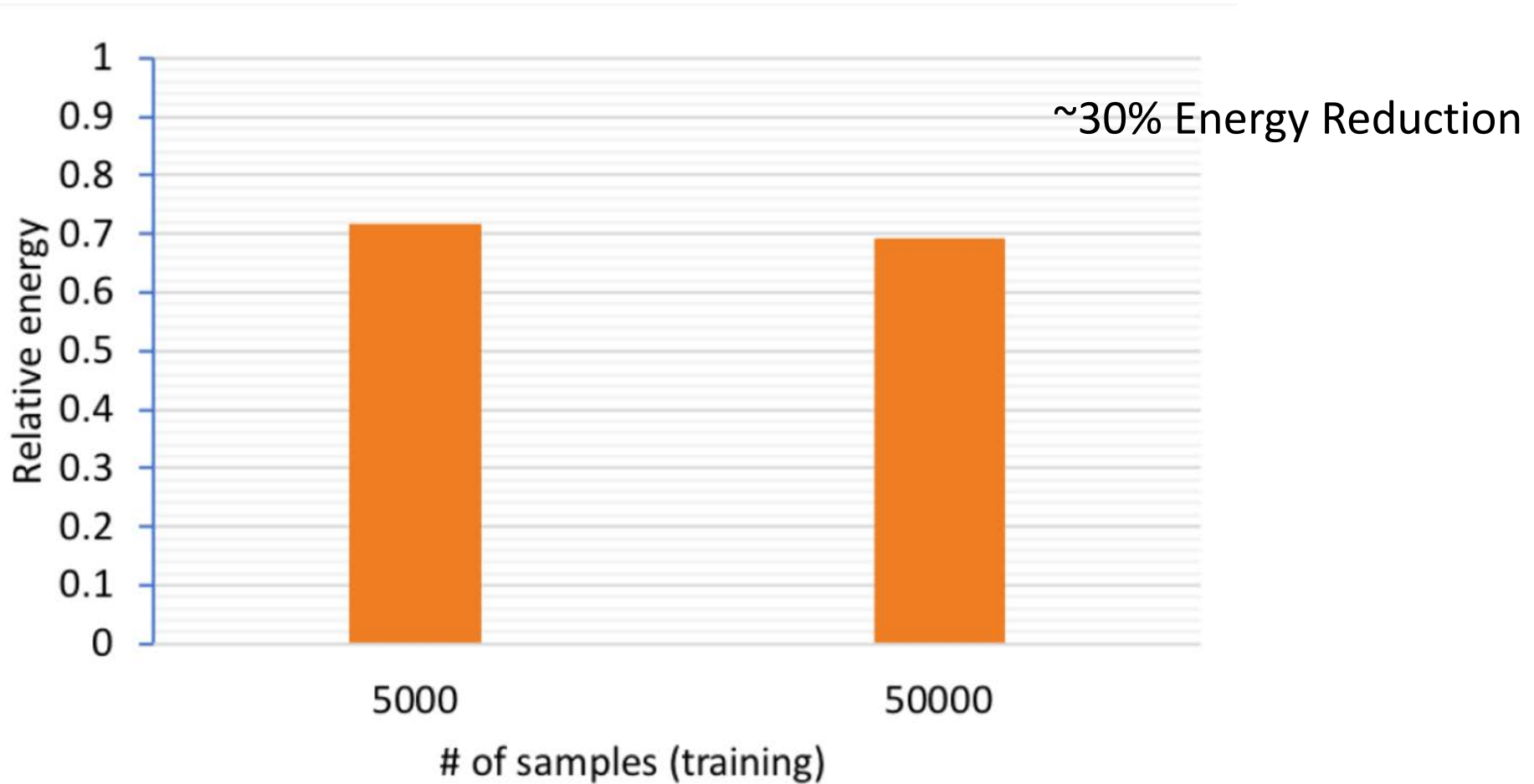
Entire Single Precision **breaks** KRLS

Identical Accuracy between double precision and transprecision KRLS

Execution Time Reduction by Transprecision KRLS



Energy Reduction by Transprecision KRLS



Conclusions

Transprecision Computing:

Minimising Execution Time without accuracy loss and HW resource increment

Transprecision Techniques brought Significant Speedups and Energy Savings to Linear Solvers and KRLS without Accuracy Loss compared to Full precision arithmetic

Transprecision Computing can be a crucial paradigm for ManyCore in order to achieves both Speedups and Energy Savings.

Acknowledgements



: Energy Efficient Transprecision Techniques
for Linear Solver Jun. 2018 – May 2020



This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 798209 (Entrans).



: Open Transprecision Computing

Jan. 2017 – Dec. 2020

This project has received funding by the European Commission Horizon 2020 research and innovation programme under grant agreement No. 732631(OPRECOMP)



Thank you very much

Any Questions?

