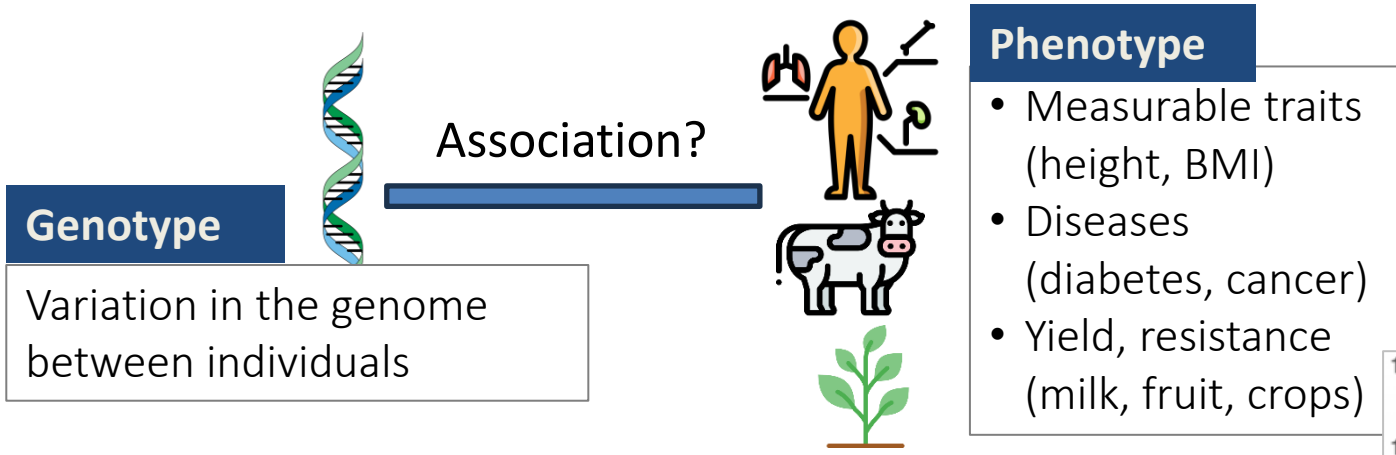


Resource-Efficient Approach for Large-Scale Genome-Wide Association Studies

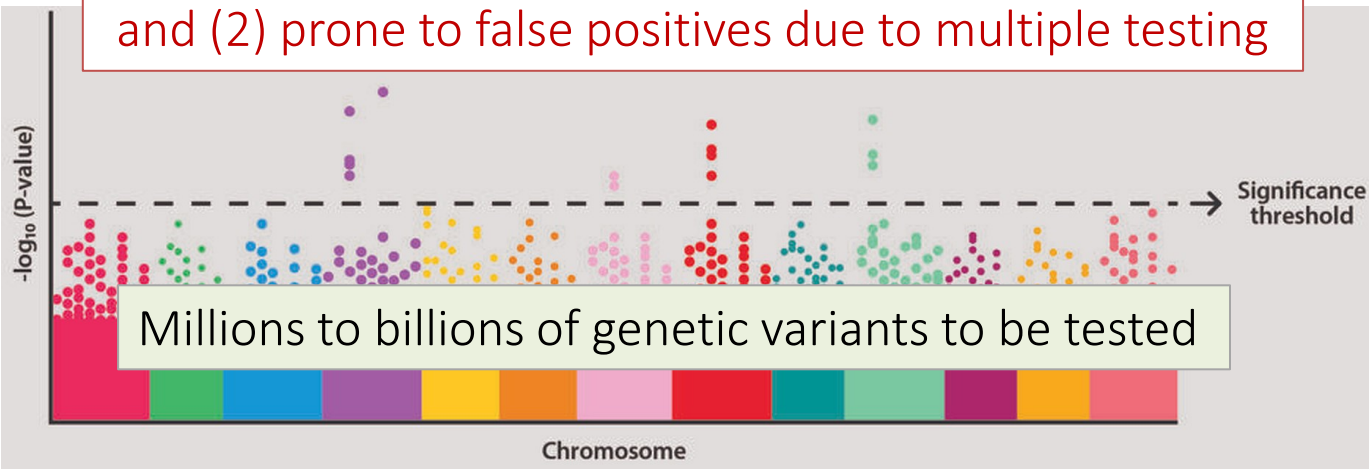
Mulya Agung

Institut Teknologi Sains Bandung

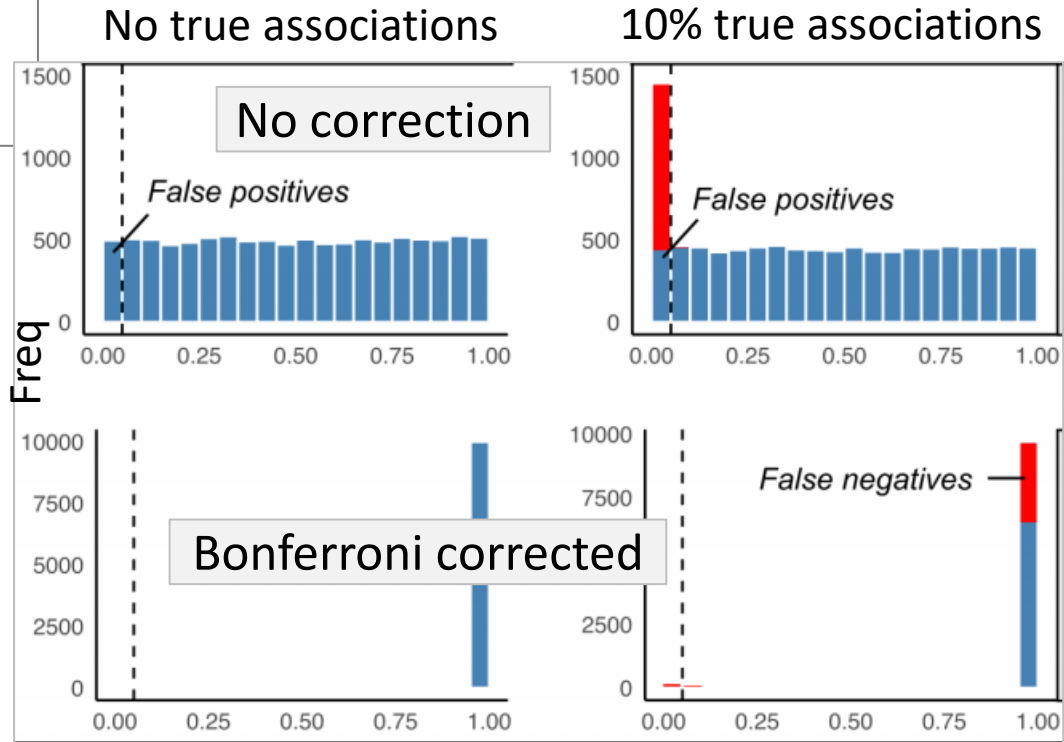
- The discovery of associations between genetic variations and traits
- More than 90,000 GWAS have been performed as of 2024 (Harris et al., 2025)



Association testing is (1) computationally expensive and (2) prone to false positives due to multiple testing

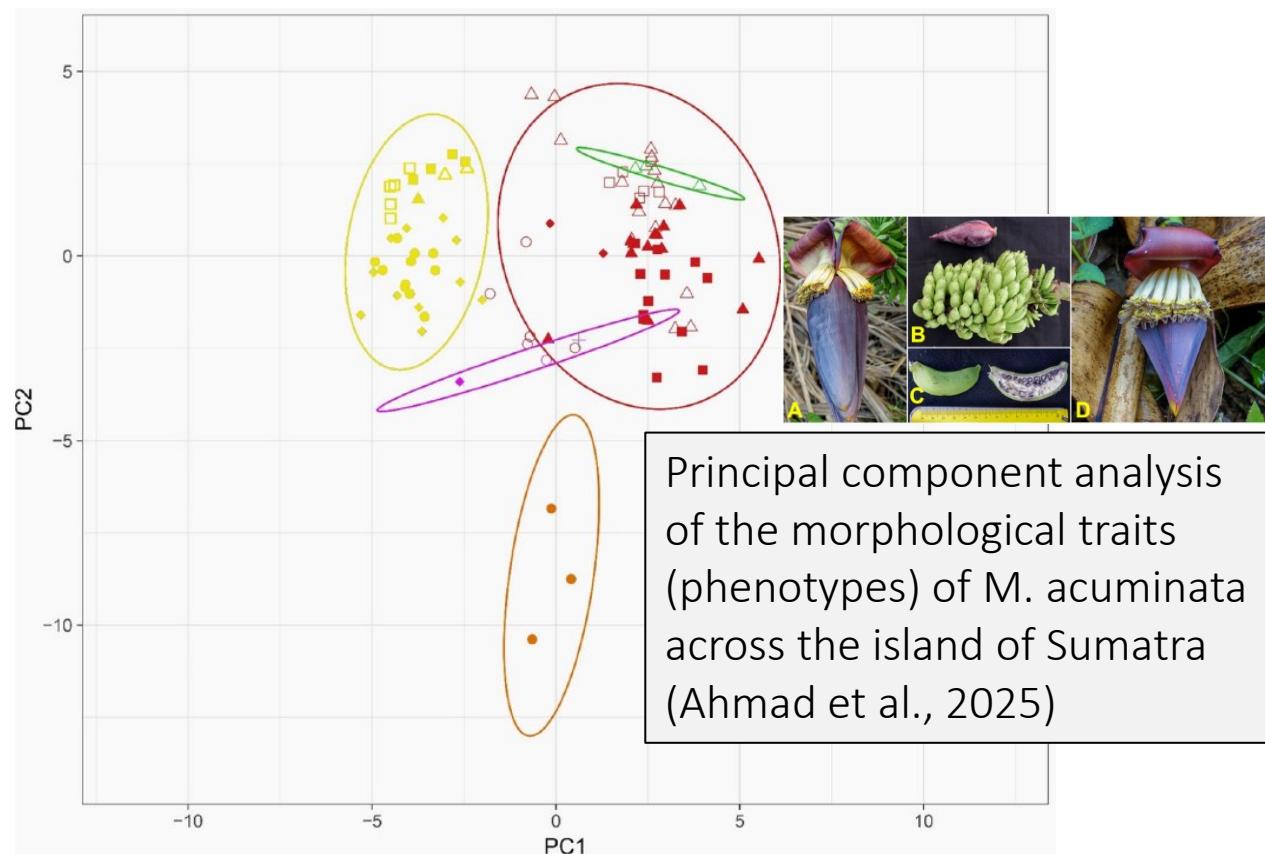


Multiple testing correction is needed to maintain false positives



- Bananas are an essential fruit worldwide due to their rich nutrients. Southeast Asia is a major part of the origin of bananas (*Musa acuminata*)
- Genetic exploration of wild ancestors of banana cultivars** is important for conservation and breeding
- Banana producers look for specific morphology, fruit quality, yield, and agronomic features
- The decrease in sequencing cost allows whole-genome sequencing in bananas** (tens of millions of variants reported)
- High genetic diversity on morphological traits across wild varieties**

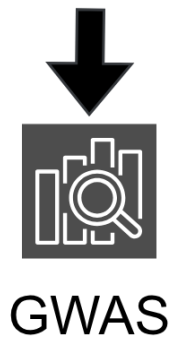
The large number of variants in whole-genome sequencing opens possibilities for discoveries, but also poses the GWAS challenges



Proposed Method: CMA: Divide-And-Conquer GWAS

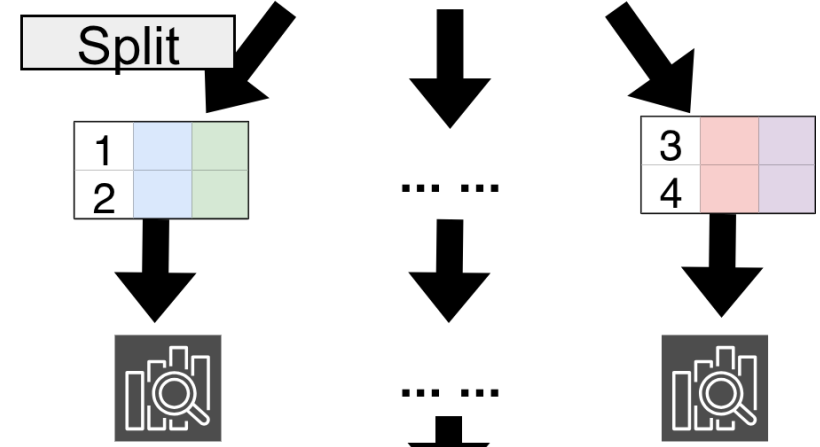
Traditional GWAS uses whole samples and variants

Sample/ accession	ID	Variants
1	1	[blue, green, orange, yellow, red, purple]
2	2	[blue, green, orange, yellow, red, purple]
3	3	[blue, green, orange, yellow, red, purple]
4	4	[blue, green, orange, yellow, red, purple]



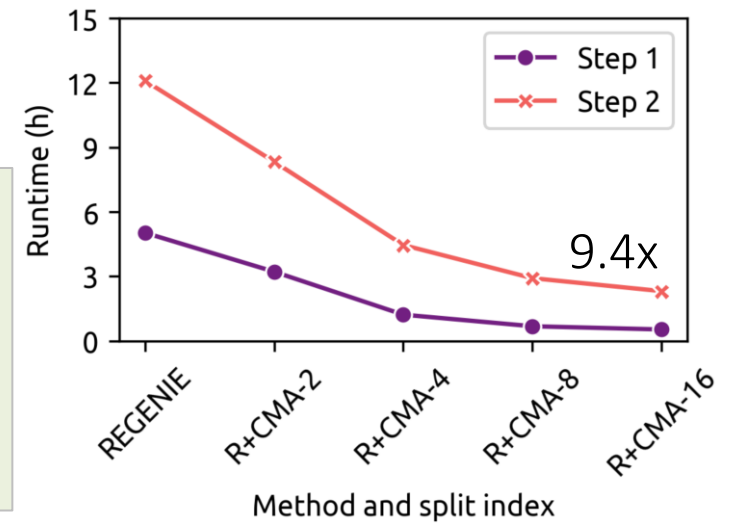
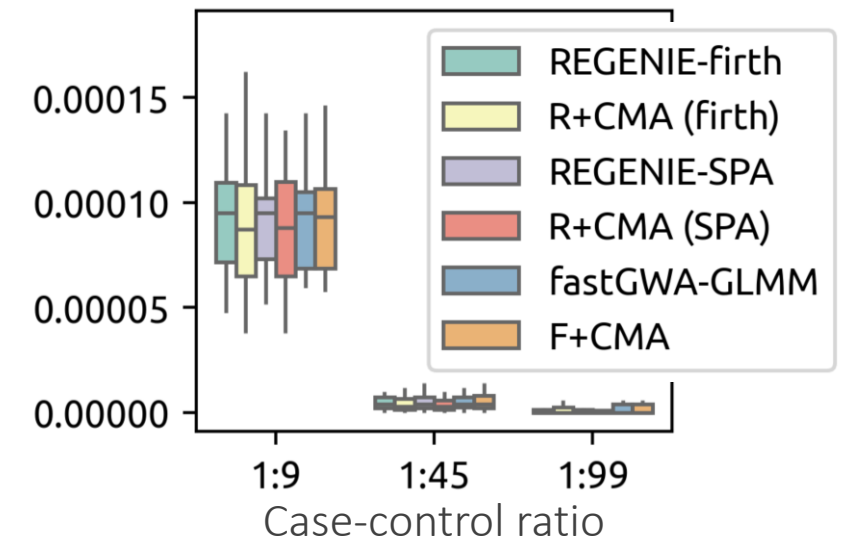
CMA is available on <https://git.ecdf.ed.ac.uk/cma>
 Publication on GENETICS: <https://doi.org/10.1093/genetics/iyaf019>

ID	Variants
1	[blue, green, orange, yellow, red, purple]
2	[blue, green, orange, yellow, red, purple]
3	[blue, green, orange, yellow, red, purple]
4	[blue, green, orange, yellow, red, purple]



CMA is much faster than the current state-of-the-art, with less memory usage

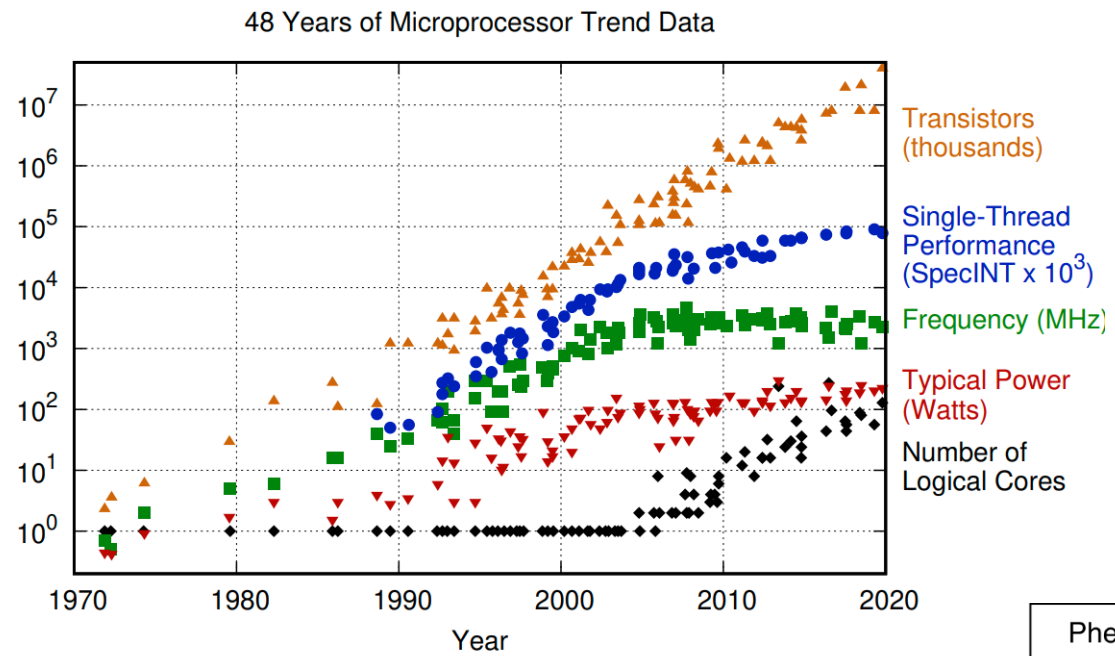
Type 1 error rates on human datasets



Proposed Method: GWA-X: GWAS permutation testing on GPUs

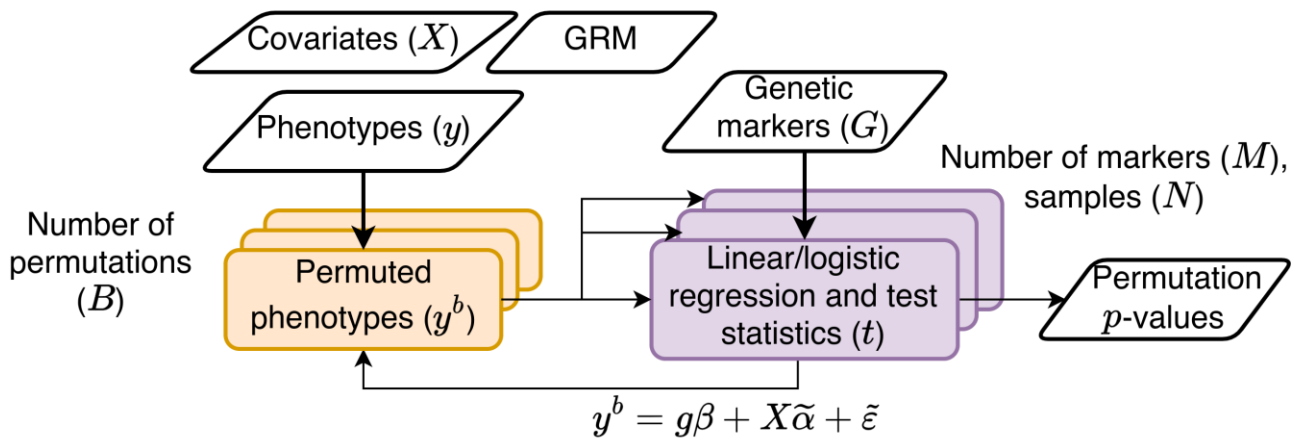
Permutation testing is the gold standard for multiple testing correction to reduce false positive rates

GWAS permutation testing is often impractical due to the large genetic data

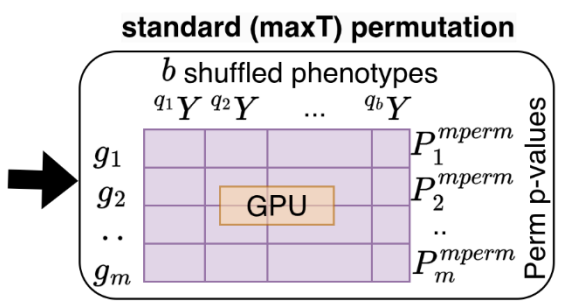
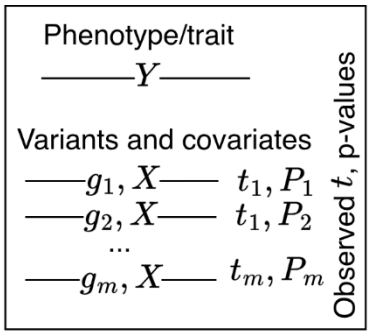


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2019 by K. Rupp

GWA-X is available on
<https://git.ecdf.ed.ac.uk/magung/gwa-x>
 Publication on bioRxiv:
<https://doi.org/10.1101/2024.09.15.613119>



The **number of regressions** increases from (M) to ($M \cdot B$)
 The complexity of matrix multiplication and inverse: $O(N^3 \cdot M \cdot B)$



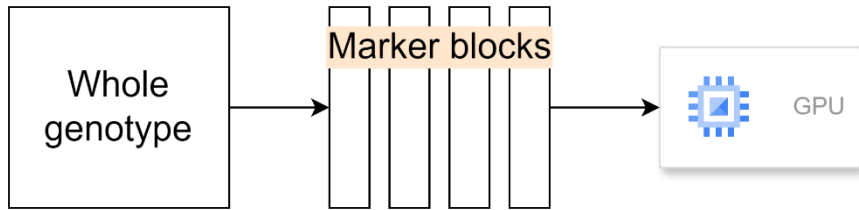
$$P_j^{mperm} = \frac{R_j}{b}, R_j \sim Binom(b, P)$$

R_j is the number of success of j th variant (${}^q t_j \geq t_j$), b is the number of permutations.

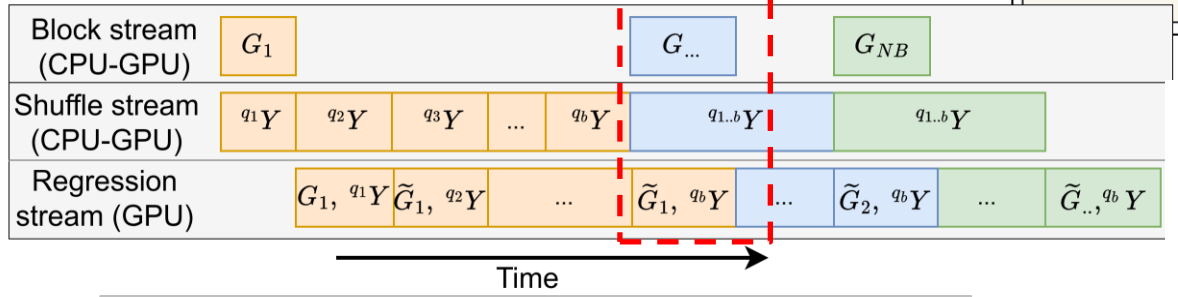
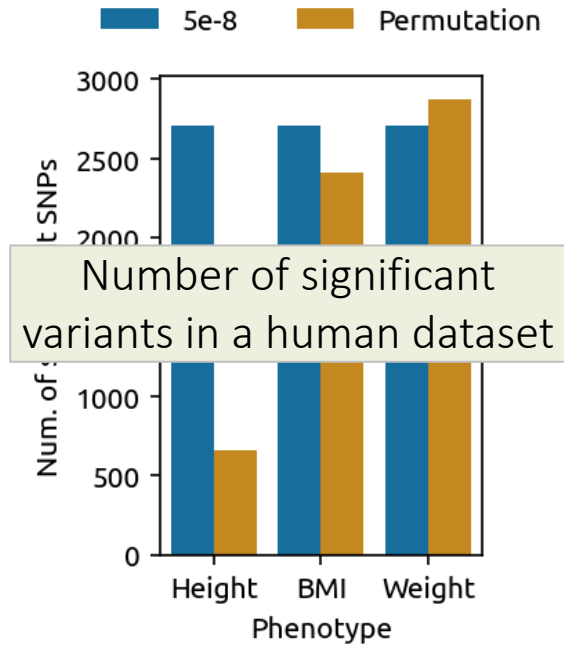
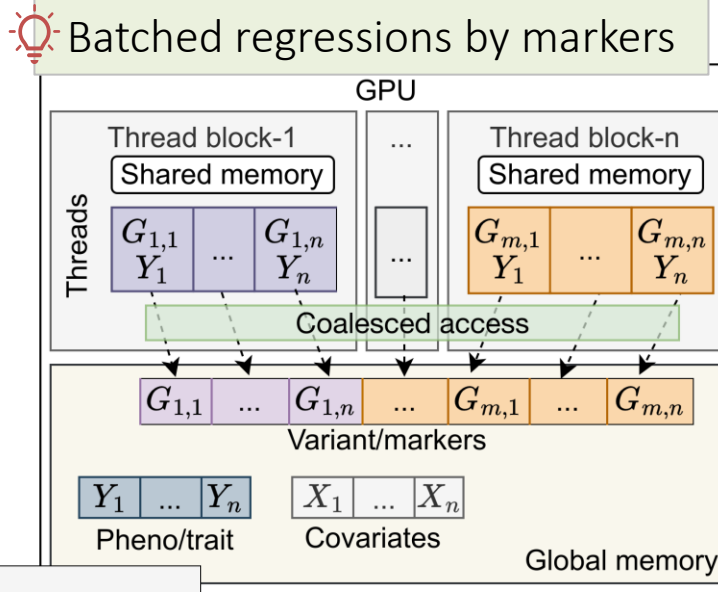
$$P_j^{aperm} = \begin{cases} \frac{r}{B_j}, & R_j < r \end{cases}$$

Proposed Method: GPU Optimization and Performance

Genotype (G) with m markers/variants, n samples/accessions

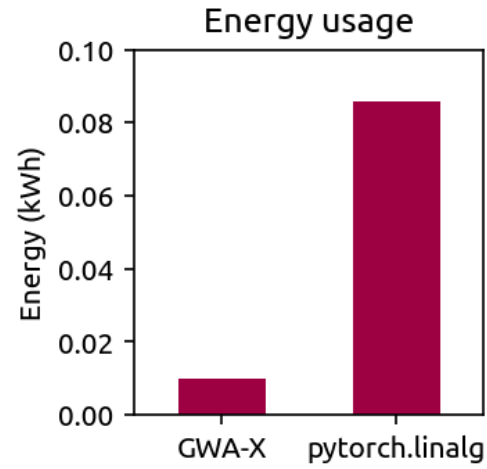
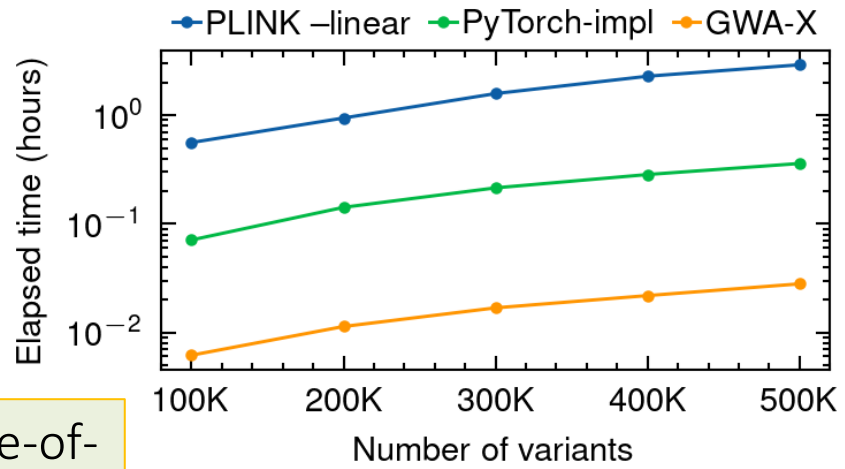


Genotype partitioning to handle a large genotype



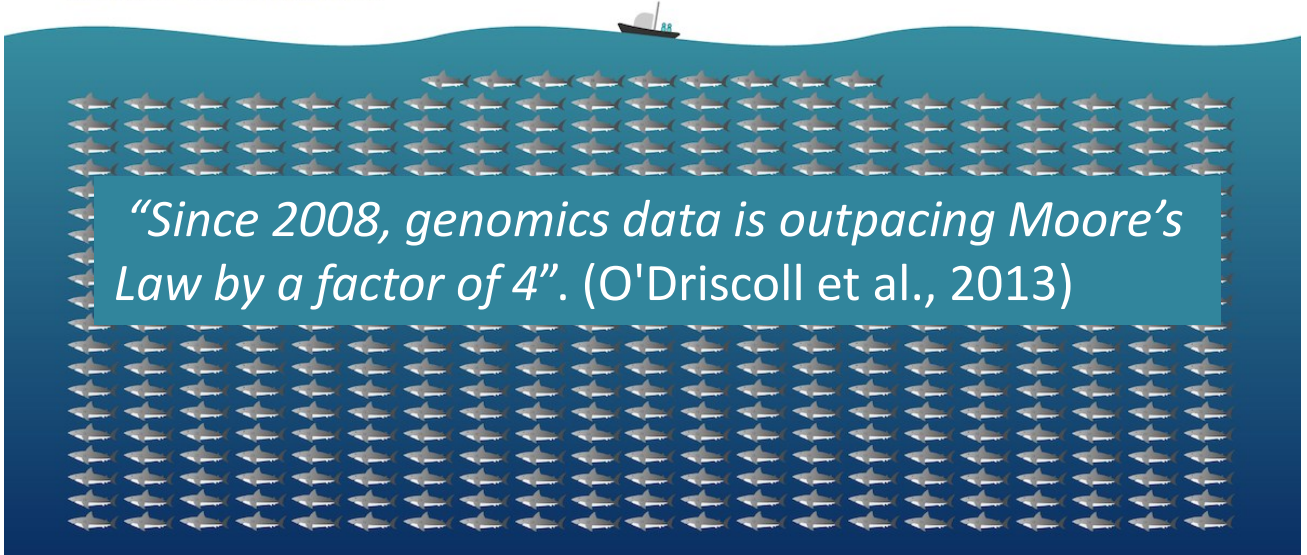
Overlapping data transfers with kernel computation to reduce data transfer delays

GWA-X is 100x times faster than the current state-of-the-art on CPUs (88.4% less energy usage)



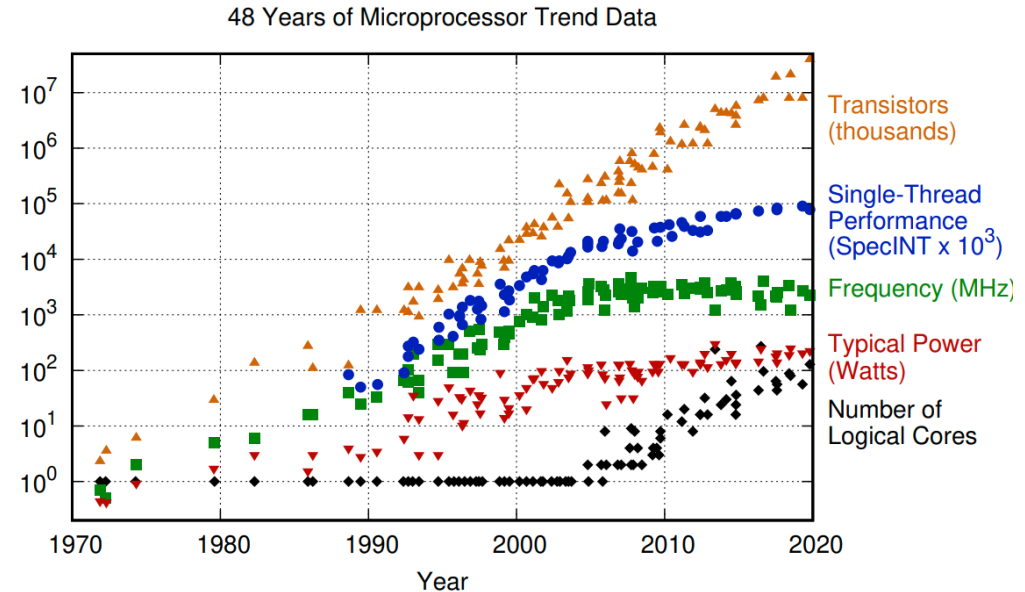
How big is 40 exabytes?

Genomics projects will generate 40 exabytes of data in the next decade.
 Each shark = 100,000,000 GB of data



“Since 2008, genomics data is outpacing Moore’s Law by a factor of 4”. (O’Driscoll et al., 2013)

Computer architecture shifts to using accelerators, e.g., GPUs



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and G. Batten. New plot and data collected for 2010-2019 by K. Rupp.

- The growth of the genomics data size scale requires a resource-efficient method scalable to recent computer architectures
- We anticipate that further advances in developing efficient software in GWAS will adopt our approach to mitigate the challenges associated with large data size

Impact: Democratize the Use of ICT Systems for Genomic Research

“Through the use of cloud computing, data commons can support large-scale data, but this also creates sustainability challenges, due to the cost of large-scale storage and compute”. (Grossman, 2019)

Databases

1982–present



- Data repository
- Researchers download data

2010–2020



- Supports large datasets and data-intensive computing with cloud computing
- Researchers can analyze data with their own virtual workspaces and applications and collaborate with other researchers with collaborative workspaces (data do not have to be downloaded to be analyzed)

Data commons

2014–2024



- Supports large datasets and data-intensive computing with cloud computing
- Workspaces
- Common data models
- Core data services
- Data and commons governance
- Harmonized data
- Data sharing
- Reproducible research

Trends in Genetics

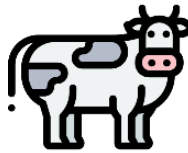
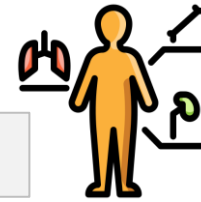
Resource-efficient GWAS will promote the sustainability goals in ICT and benefit the research community, especially researchers with limited resources and budgets

Genotype

Variation in the genome between individuals



GWAS: CMA, GWA-X



Phenotype

- Measurable traits (height, blood pressure)
- Diseases (diabetes, cancer)
- Yield, resistance (milk, fruit, crops)

- We anticipate collaboration between researchers from multiple disciplines conducting **genomic research in humans, animals, and plants**
- Our GWAS project on wild bananas is one example of the collaboration
 - **Identifying the genetic variants controlling agronomic traits** in plants is important for developing new varieties with desirable traits, e.g., higher yield and quality

Research output: GWAS Pipeline for Bananas

Public dataset (genotype-phenotype)

- Trait focus: *Fusarium* resistance, morphology, fruit quality, and yield

Collab: BRIN's Botany lab

1. Data collection (*Musa acu.*)

2. Genotyping

3. Phenotyping

4. Quality control

Population structure analysis

LD

5. Association testing (CMA, GWA-X)

GWAS summary statistics

e. Post-GWAS analysis

Fine-mapping

Meta-analysis

Heritability (h^2)

Category	Trait
Morphology	Pseudostem height
	Single leaf width
	Time from flowering to harvest
Fruit quality	pH
	Total soluble solids
	Titrateable acidity
Yield	Hands weight
	Number of hands in a bunch
	Fruit length
	Pulp percentage
	Pulp dry weight percentage
	Peel thickness

Genetic basis of *M. acuminata*

Linkage and kinship reference across varieties

- GWAS is a powerful tool for unravelling genetic variants associated with phenotypes or traits of interest, such as morphological traits, diseases, and agronomic traits
- The emerging scale of genomic data size opens possibilities for discoveries, but also poses a big data challenge
- We present a resource-efficient approach, consisting of CMA and GWA-X, to mitigate the GWAS challenge associated with the large-scale data size, including
 - A divide-and-conquer method for partitioning a large GWAS workload into workloads of smaller subsets to be run on smaller-scale computers
 - A GWAS permutation testing approach on GPUs
- One of our ongoing projects is to adopt our approach for GWAS on wild bananas (*M. acuminata*), an essential fruit in the ASEAN region, aimed at investigating genetic variants associated with disease resistance and better yield
- We anticipate that further development of GWAS tools will be based on our approach for the large datasets of humans, plants, and animals

Thank you!

CMA and GWA-X are open source
You are welcome to contribute!



CMA



GWA-X

Thanks to our collaborators:

CMA & GWA-X

- Tenesa group (University of Edinburgh, UK)
- Dr. Hyojung Paik (KISTI, South Korea)

Genomic diversity of Indonesian wild bananas

- Fajaruddin Ahmad (Research Center of Applied Botany, BRIN)

