

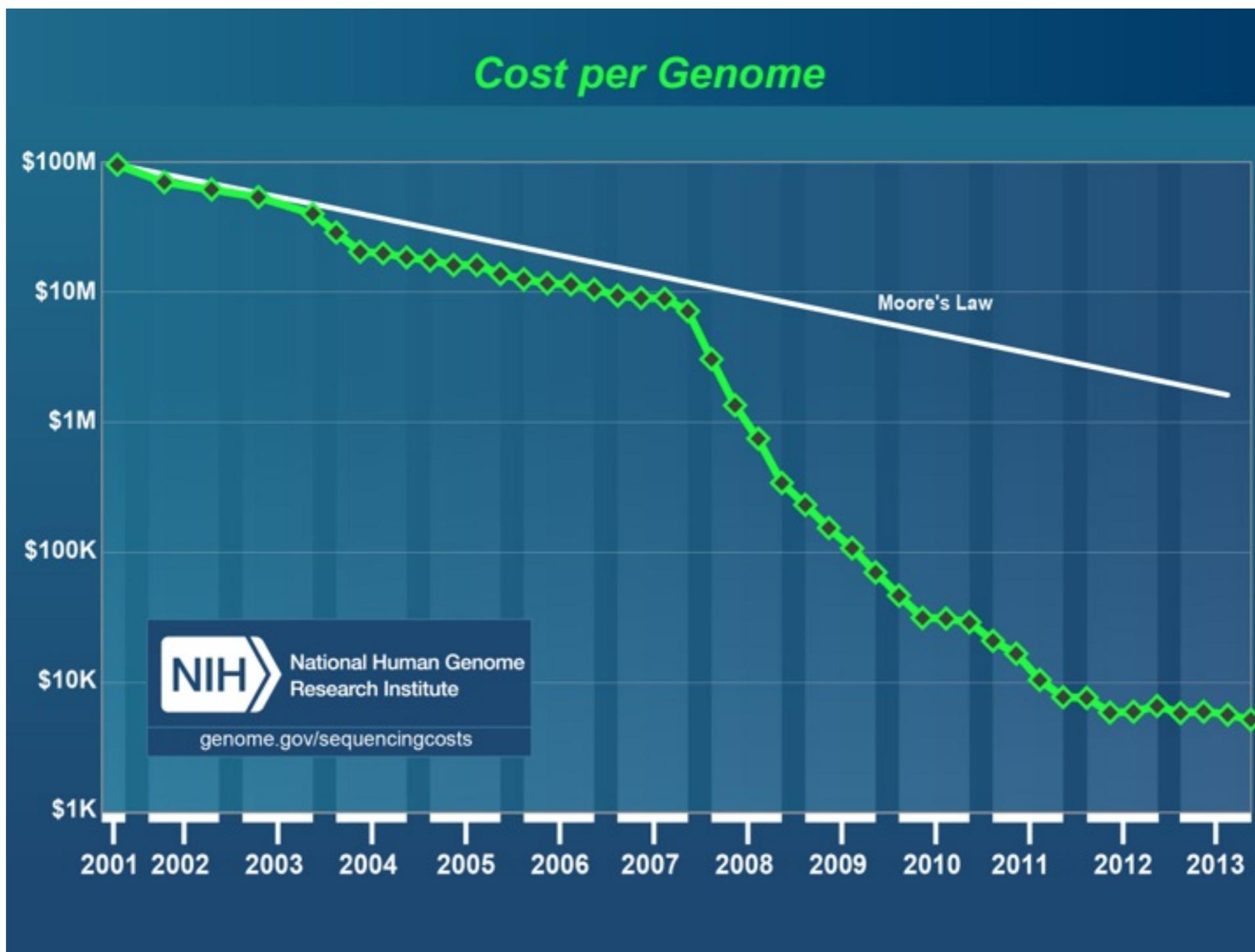
Accelerating the DNA sequencing pipeline with GPUs

Mauricio Carneiro

carneiro@broadinstitute.org

Group Lead, Computational Technology Development
Broad Institute of MIT and Harvard

vastly outpacing moor's law



Genomics Platform in 2013

50
HiSeqs

10
MiSeqs

2
NextSeqs

14
HiSeq X

6.5
Pb of data

427
projects

180
people

2.1
Tb/day



Genomics Platform in 2013

44,130
exomes

2,484
exome express

2,247
genomes

2,247
assemblies

8,189
RNA

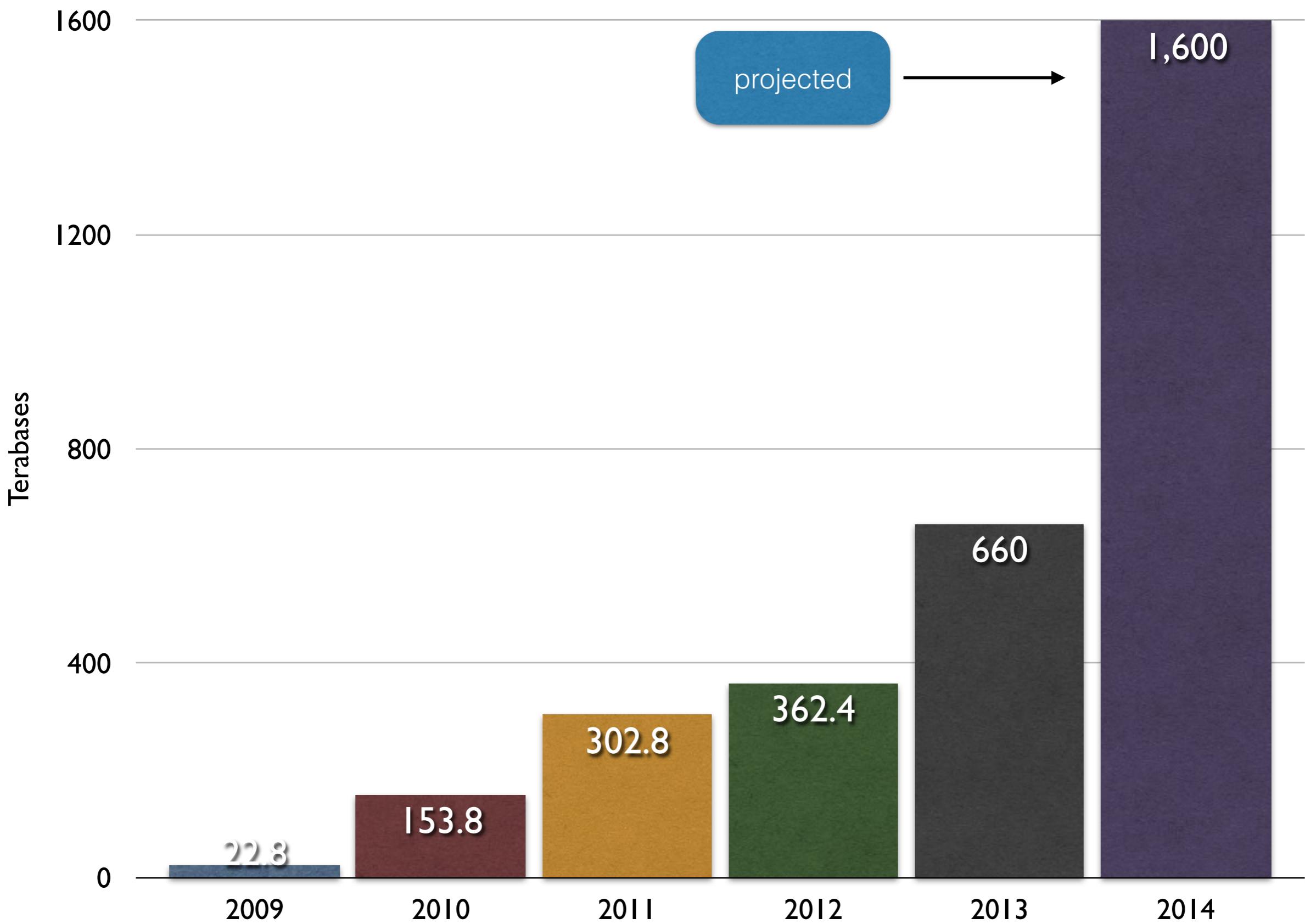
9,788
16S

47,764
arrays

228
cell lines

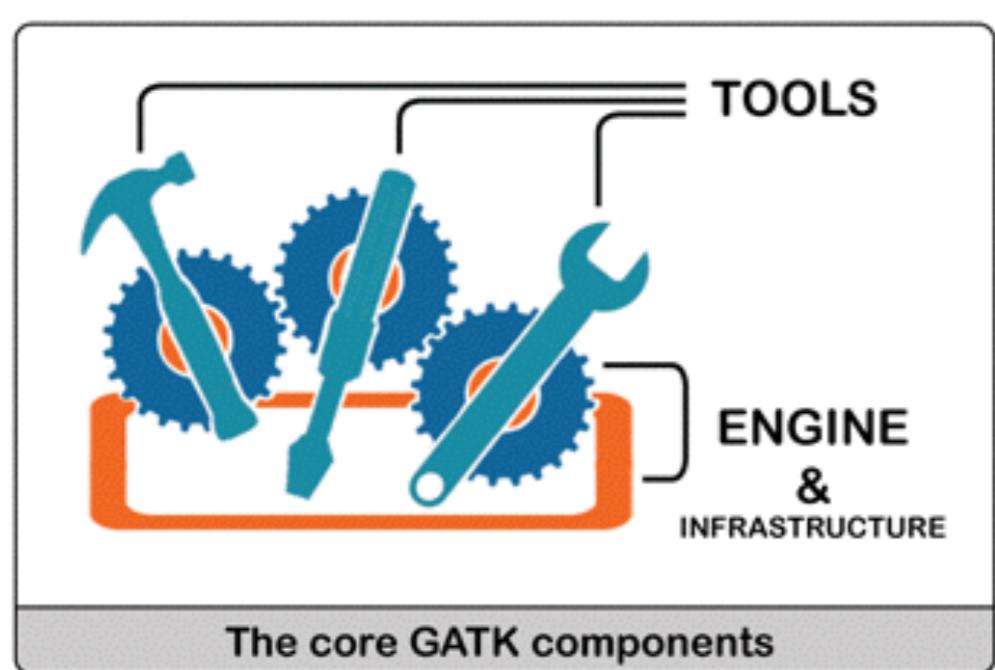


Terabases of Data Produced by Year



GATK is both a toolkit and a programming framework, enabling NGS analysis by scientists worldwide

Toolkit & framework packages

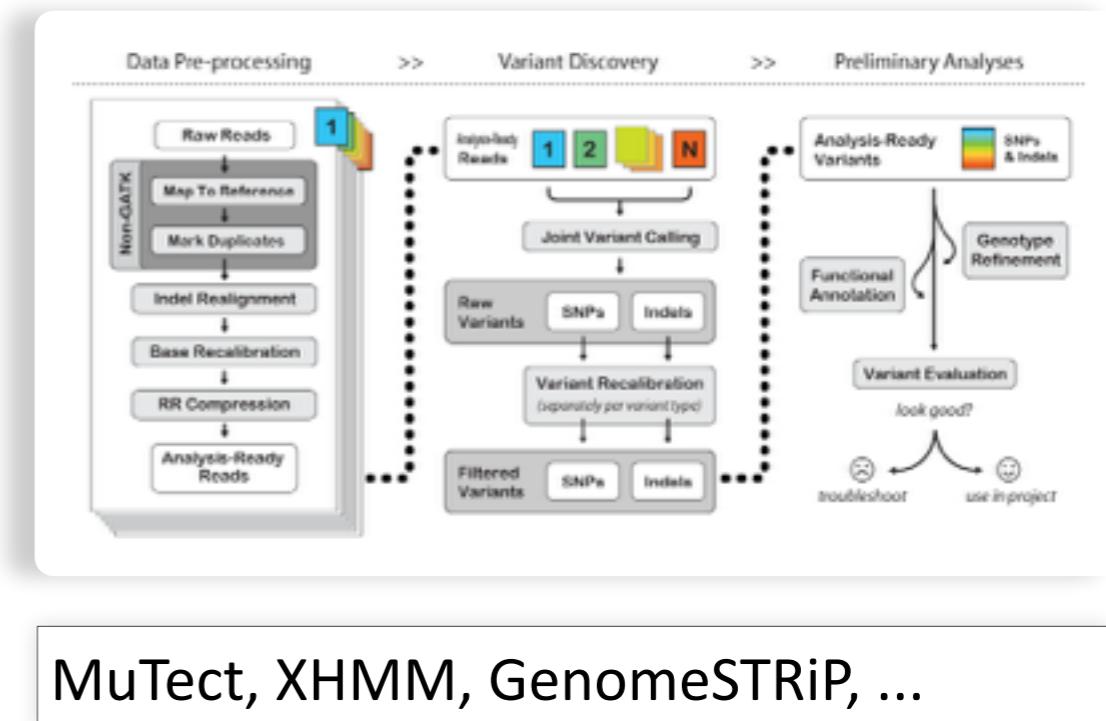


Toolkit
→

*Best practices
for variant
discovery*

Framework
→

Tools developed on top of the GATK framework by other groups



MuTect, XHMM, GenomeSTRiP, ...

Extensive online documentation & user support forum serving >10K users worldwide



<http://www.broadinstitute.org/gatk>



About

Overview of the GATK and the people behind it



Guide

Detailed documentation, guidelines and tutorials



Community

Forum for questions and announcements



Events

Materials from live and online events

Workshop series educates local and worldwide audiences

Completed:

- Dec 4-5 2012, Boston
- July 9-10 2013, Boston
- July 22-23 2013, Israel
- Oct 21-22 2013, Boston

Planned:

- March 3-5 2014, Thailand
- Oct 18-29 2014, San Diego

iTunes U Collections



BroadE: GATK
Broad Institute



Format

- Lecture series (general audience)
- Hands-on sessions (for beginners)

Portfolio of workshop modules

- GATK Best Practices for Variant Calling
- Building Analysis Pipelines with Queue
- Third-party Tools:
 - GenomeSTRiP
 - XHMM

Tutorial materials, slide decks and videos all available online through the GATK website, YouTube and iTunesU

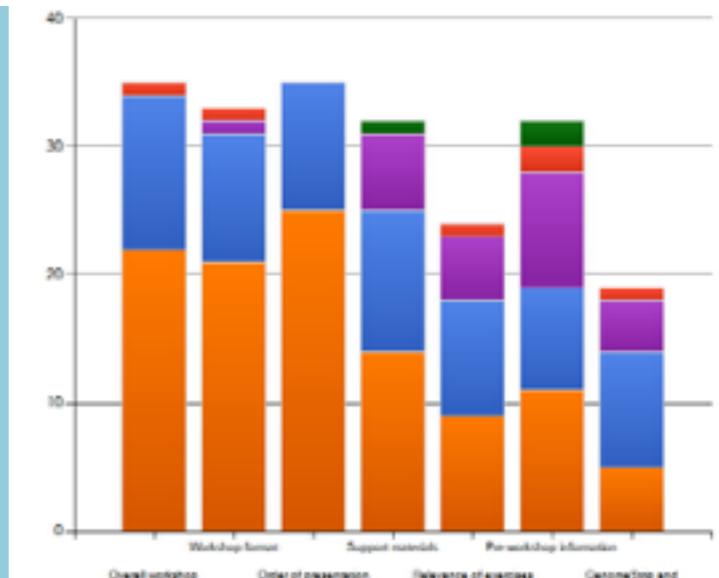
BroadE: Overview of GATK & best practices

by broadinstitute • 1 week ago • 1 view

Copyright Broad Institute, 2013. All rights reserved. The presentations below were filmed during the 2013 GATK Workshop, part of ...

NEW HD

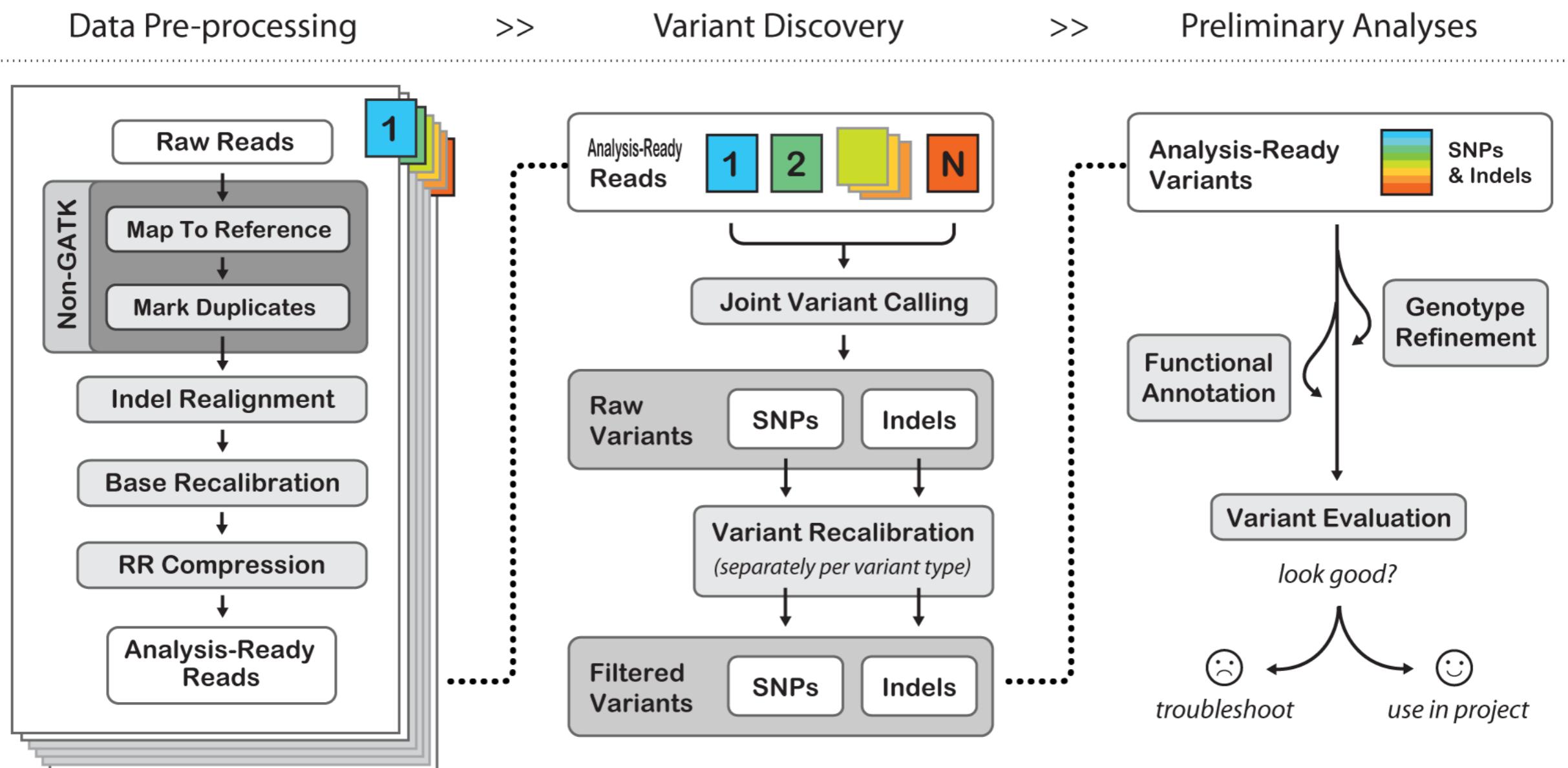
22:06



- High levels of satisfaction reported by users in polls
- Detailed feedback helps improve further iterations



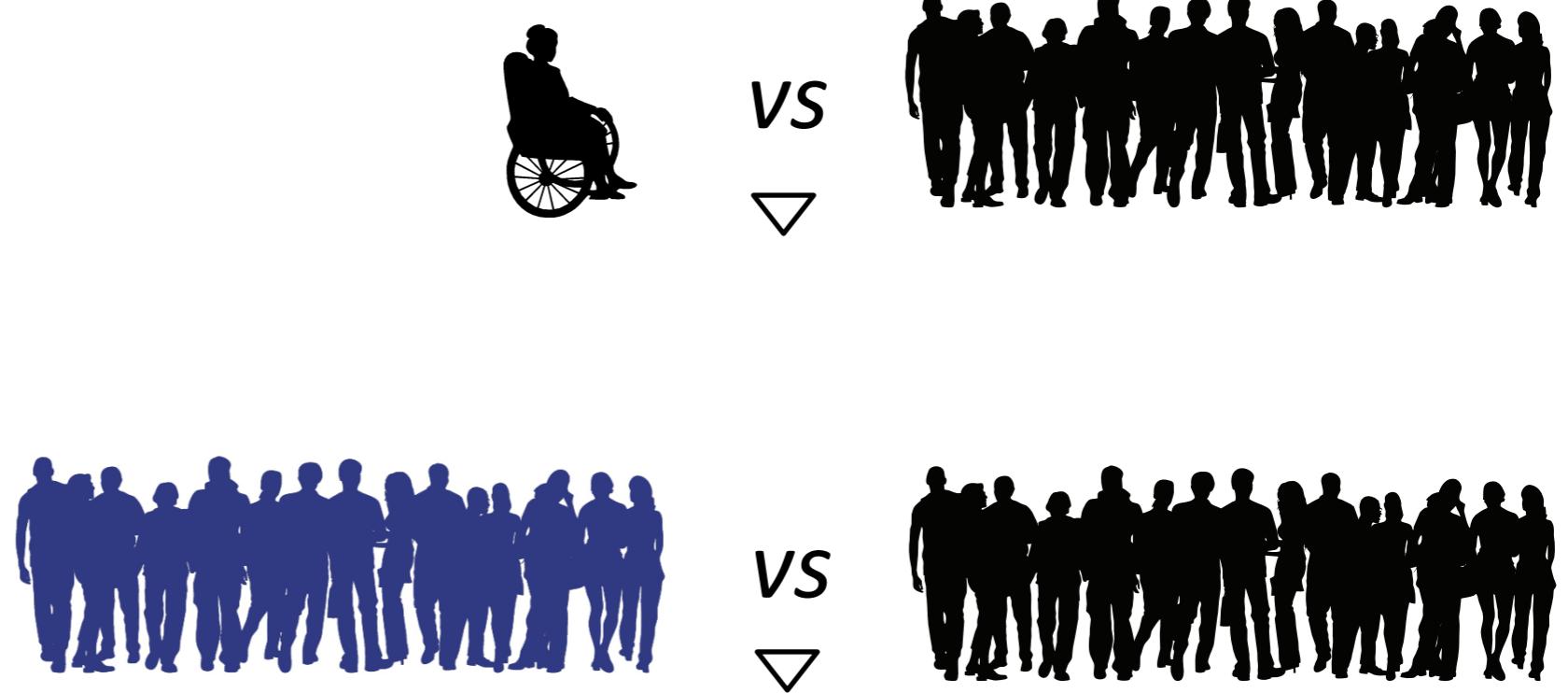
We have defined the best practices for sequencing data processing



To fully understand **one** genome we need **tens of thousands** of genomes

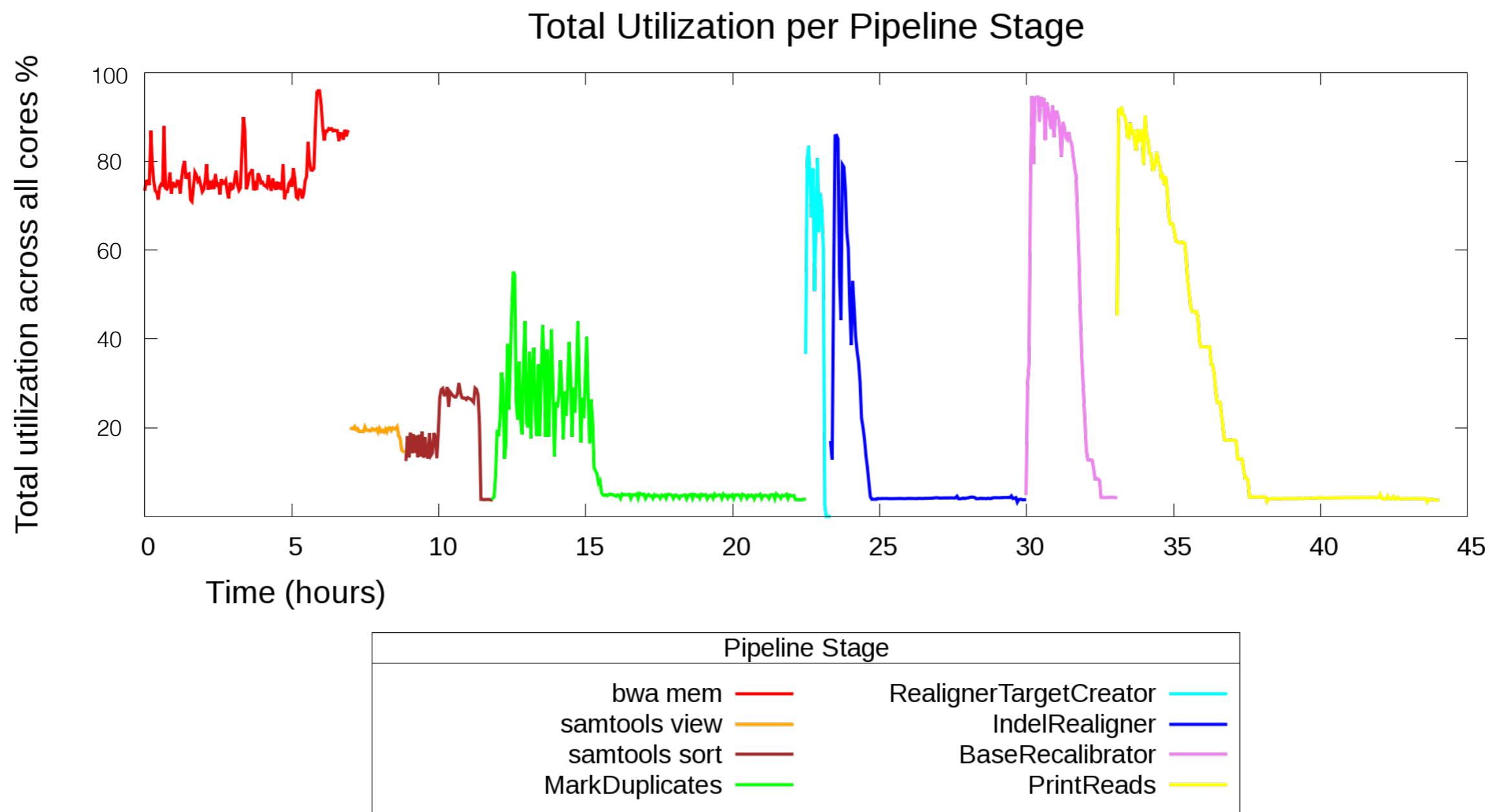
Rare Variant
Association Study
(RVAS)

Common Variant
Association Study
(CVAS)



Technical challenge
all samples need to be consistently processed

Processing times on a whole genome sample



Future challenges to scale up the pre-processing pipeline

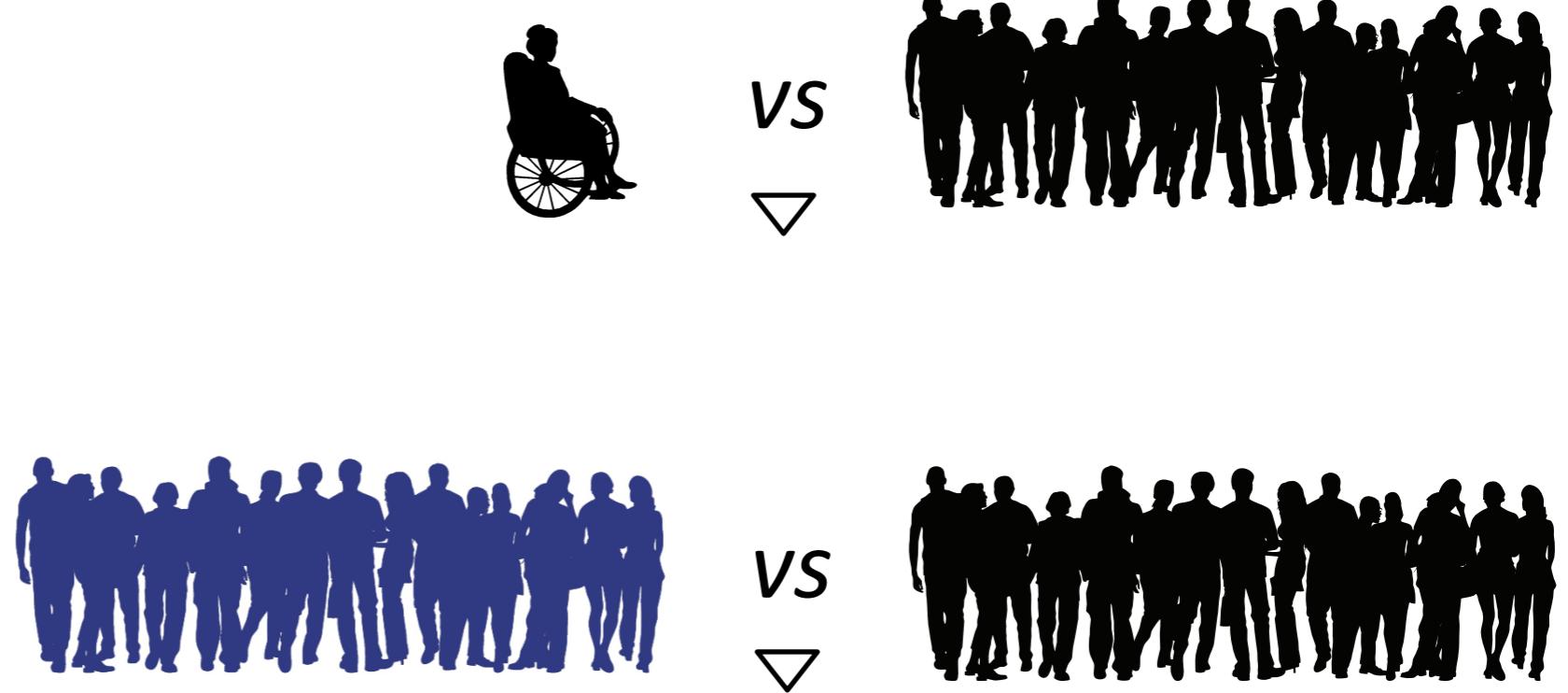
- Eliminate disk read/writing in between pipeline steps
- Redesign algorithms with performance in mind (e.g. Duplicate Marking)
- Make use of heterogeneous compute alternatives (e.g. SSE/AVX, GPU, ...)
- Reduce time spent doing unnecessary calculations (e.g. Base Recalibration on *good* data)
- Eliminate Indel Realignment
- New methods (entirely) that can replace long running processes (e.g. different statistical approach for Base Recalibration?)

step	threads	time
BWA	24	7
samtools view	1	2
sort + index	1	3
MarkDuplicates	1	11
RealignTargets	24	1
IndelRealigner	24	6.5
BaseRecalibrator	24	1.3
PrintReads + index	24	12.3
Total		44

To fully understand **one** genome we need **tens of thousands** of genomes

Rare Variant
Association Study
(RVAS)

Common Variant
Association Study
(CVAS)



Technical challenge
all samples must be *jointly called*

The ideal database for RVAS and CVAS studies would be a complete matrix

A diagram illustrating a variant database matrix. A vertical blue arrow on the left points downwards, labeled $\sim 3M$ variants. A horizontal blue arrow at the top points to the right, labeled All case and control samples.

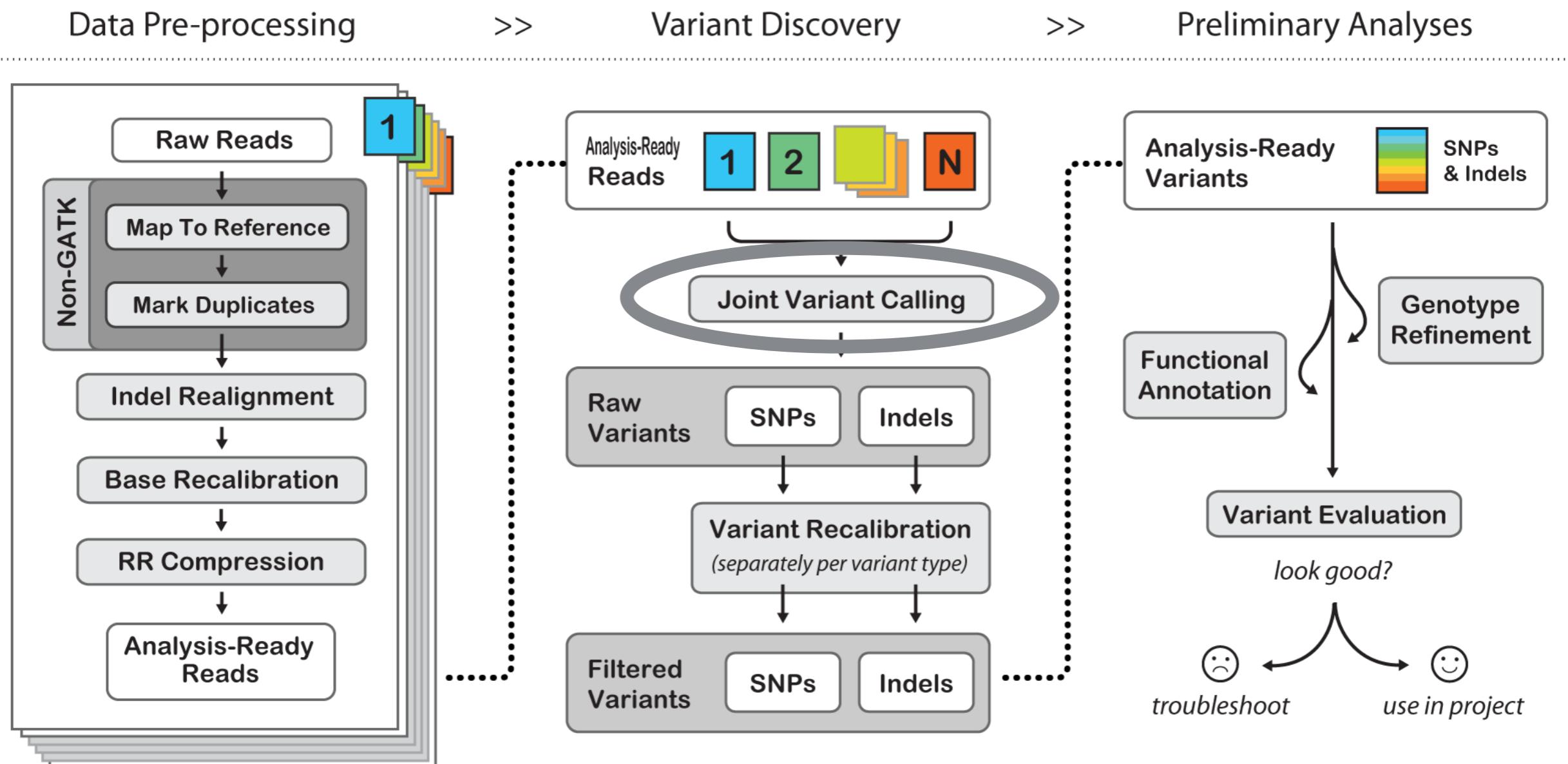
The matrix has columns for Site, Variant, Sample 1, Sample 2, ..., and Sample N. The rows are categorized by variant type: SNP, Indel, SNP, SNP, and SNP. Each row contains phred-scaled probability data for each sample.

Genotypes:
0/0 ref
0/1 het
1/1 hom-alt

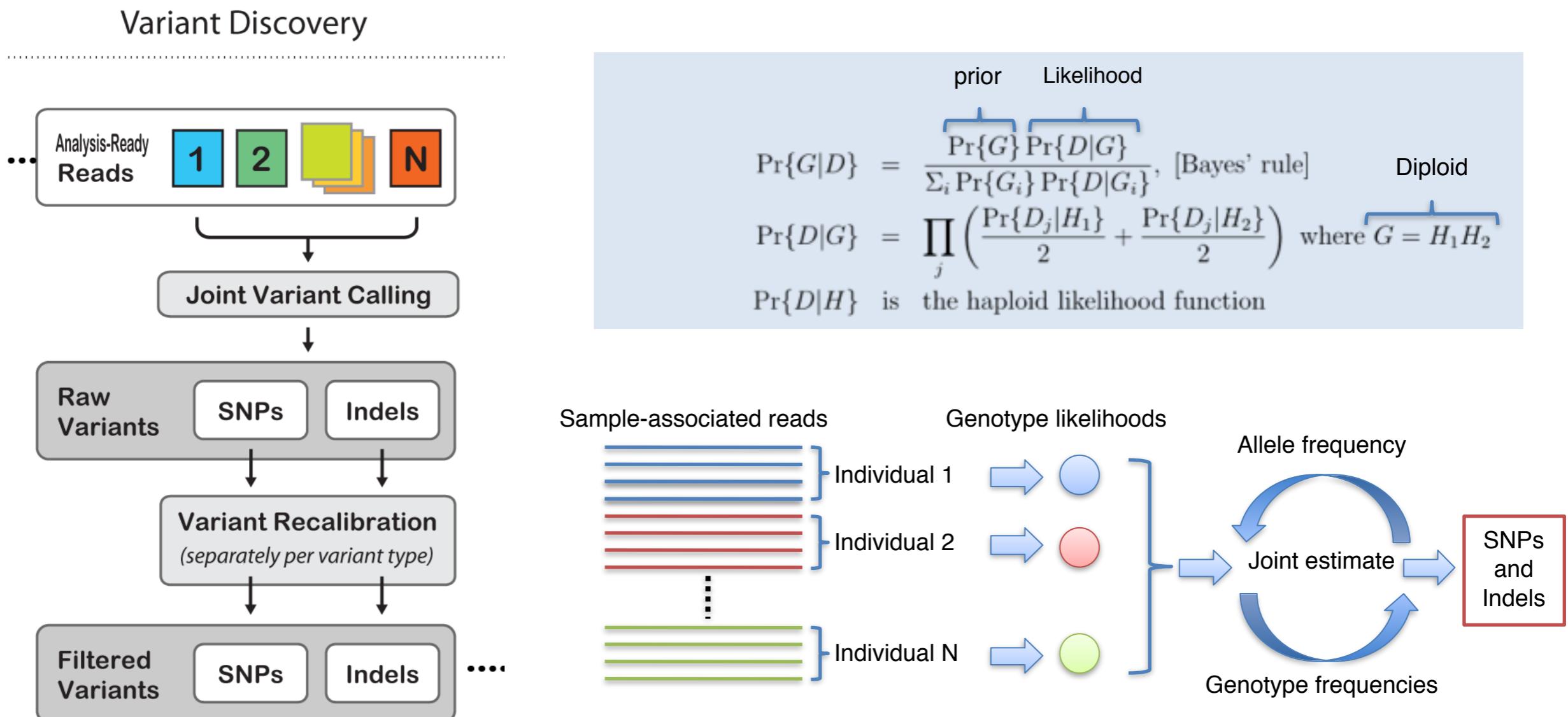
Likelihoods:
A/B/C phred-scaled probability of hom (A), het (B), hom-alt (C) genotypes given NGS data

		All case and control samples						
		Site	Variant	Sample 1	Sample 2	...	Sample N	
SNP	1:1000	A/C		<i>0/0</i> 0,10,100	<i>0/1</i> 20,0,200	...	<i>0/0</i> 0,100,255	
Indel			1:1050	T/TC	<i>0/0</i> 0,10,100	<i>0/0</i> 0,20,200	...	<i>1/0</i> 255,0,255
SNP	1:1100	T/G		<i>0/0</i> 0,10,100	<i>0/1</i> 20,0,200	...	<i>0/0</i> 0,100,255	
SNP	
	X:1234	G/T		<i>0/1</i> 10,0,100	<i>0/1</i> 20,0,200	...	<i>1/1</i> 255,100,0	

Joint calling is an important step in Variant Discovery



Variant calling is a large-scale bayesian modeling problem

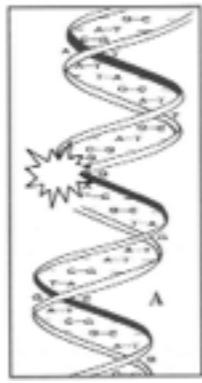


Calling 57,000 samples jointly is not easy

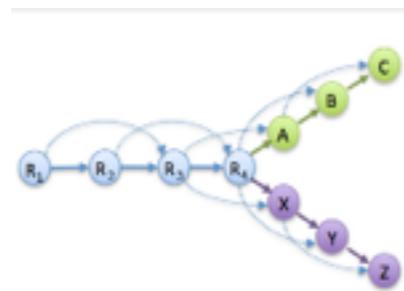
- Raw BAM files are **~1 Pb**, Reduced BAM files are **~100 Tb**
- **2,508,000 CPU hours** for data processing and **84,000** for variant calling with the *Unified Genotyper*.
- Some regions required **64 Gb RAM** (most required under 16 Gb)
- Final compressed VCFs are **2.35 Tb**

Understanding the Haplotype Caller

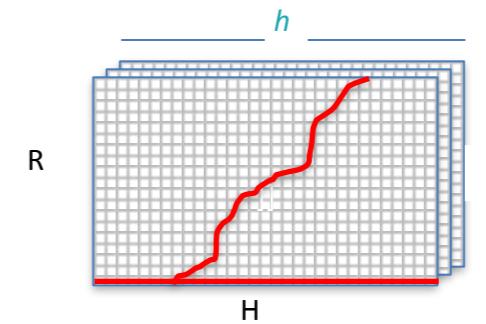
Understanding the Haplotype Caller



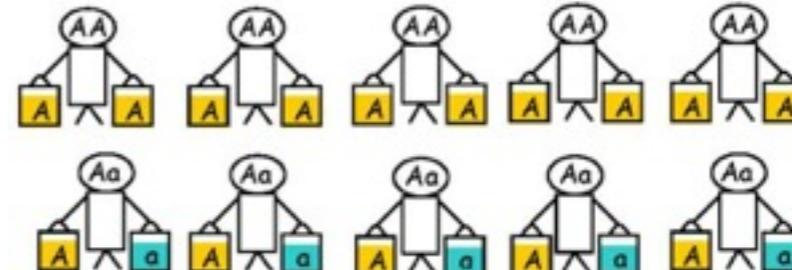
1. Active region traversal
identifies the regions that need
to be reassembled



2. Local de-novo assembly
builds the most likely
haplotypes for evaluation



3. Pair-Hmm evaluation of
all reads against all
haplotypes
(scales exponentially)



4. Genotyping
using the exact model

Pair-HMM is the biggest culprit for the low performance of the Haplotype Caller

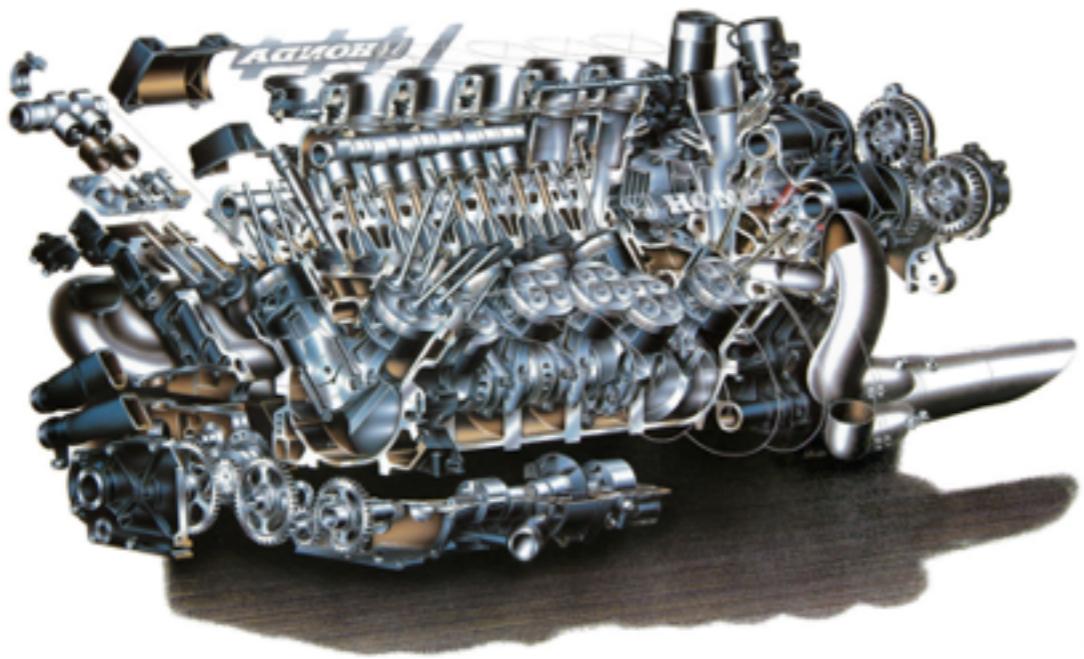
Stage	Time	Runtime %
Assembly	2,598s	13%
Pair-HMM	14,225s	70%
Traversal + Genotyping	3,379s	17%

NA12878 80xWGS performance on a single core
chr20 time: 5.6h
whole genome: 7.6 days

How can we improve the Haplotype Caller performance

- **Distributed parallelism: Queue/MapReduce**
- Alternative way to calculate likelihoods
- Heterogeneous parallel compute

The GATK is composed of two distinct parts



Engine

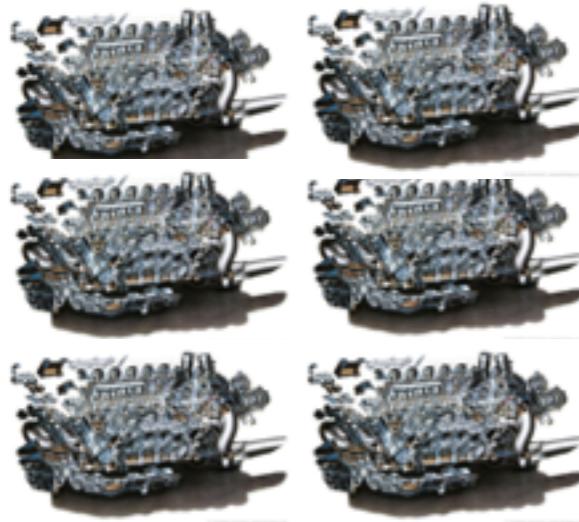
Takes care of the input and output. Preprocess and organizes a traversal system for the walkers



Walkers

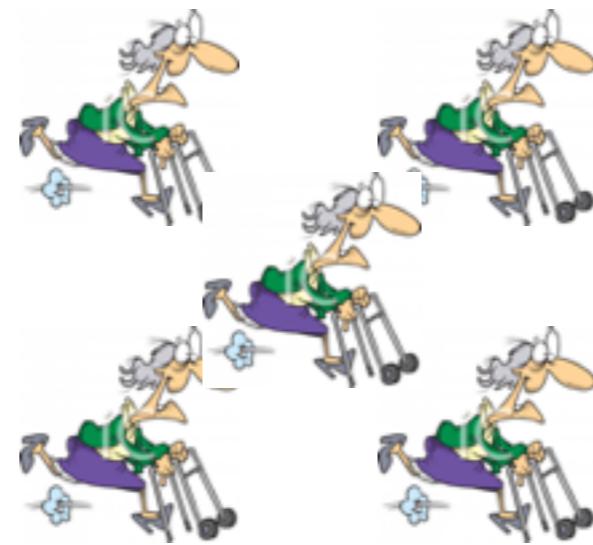
Sees the genome in an organized fashion and applies an algorithm to it

The three ways to parallelize the GATK



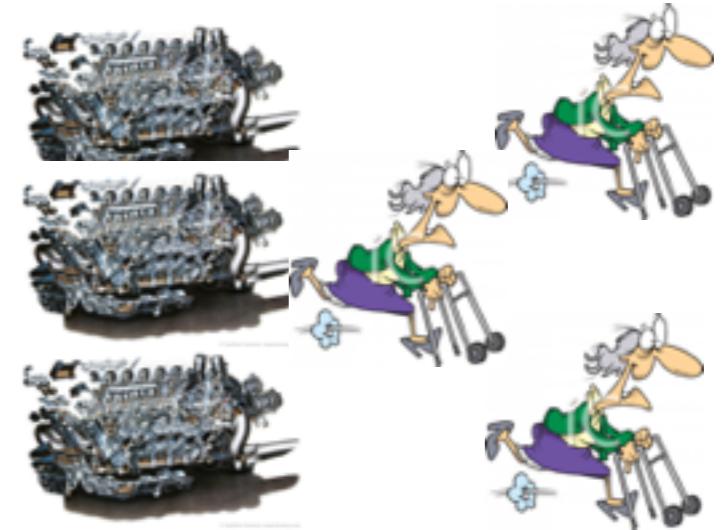
Parallelizing at the **engine** level to process different parts of the genome at the same time.

-nt



Parallelizing at the **walker** level to speed up the processing of each individual region of the genome.

-nct



Spawn many **instances** of the GATK to work on separate (arbitrary) parts of the genome at the same time.

Queue/MapReduce

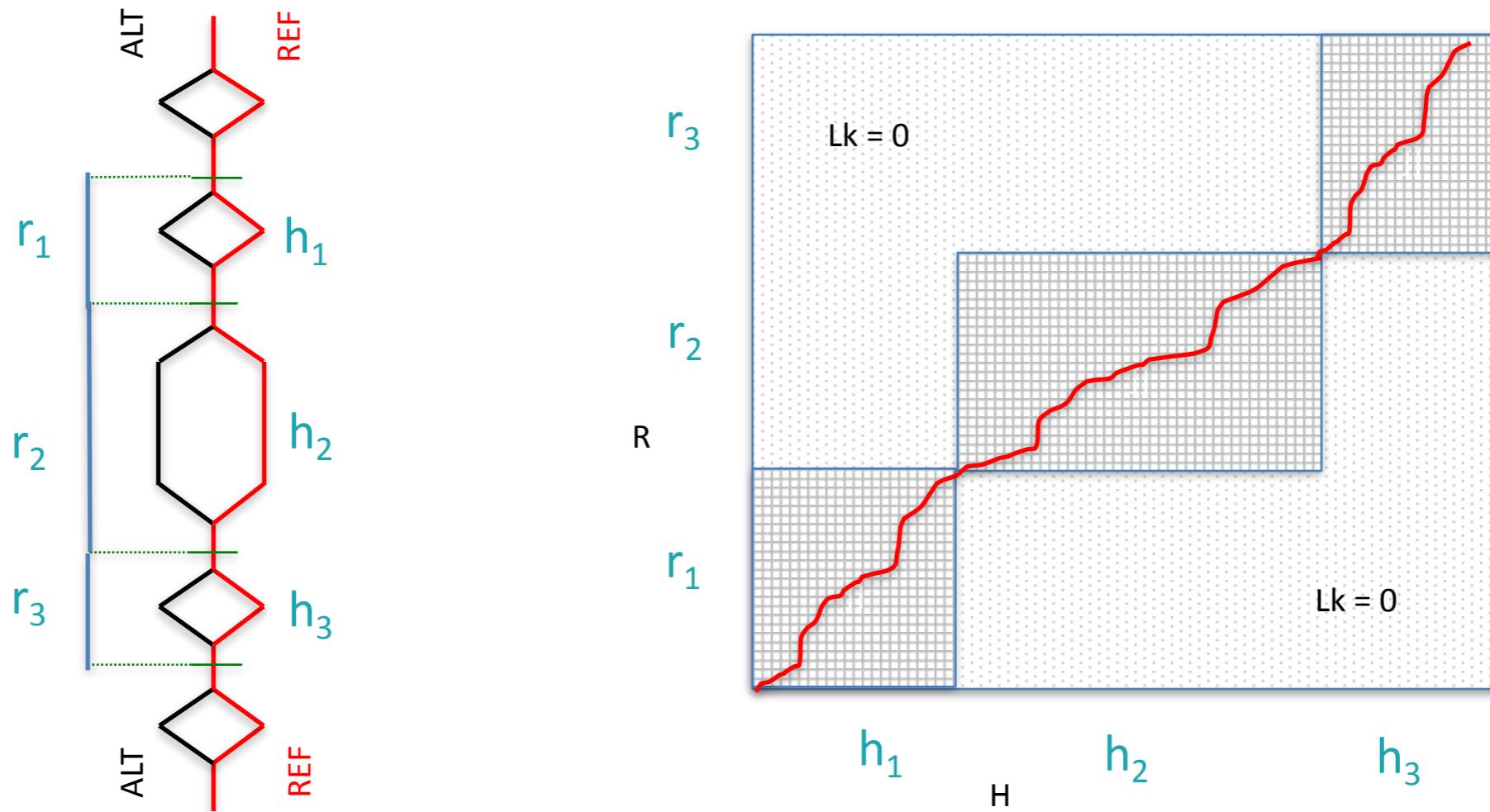
this is not really a solution, but a last resort that we use routinely

How can we improve the Haplotype Caller performance

- Distributed parallelism: Queue/MapReduce
- **Alternative way to calculate likelihoods**
- Heterogeneous parallel compute

Calculating genotype likelihoods straight from the assembly graph reduces pair-HMM usage

Mapping each read to the haplotype assembly graph we can constrain the underlying pair-HMM to avoid quasi-zero likelihood unrealistic alignments.

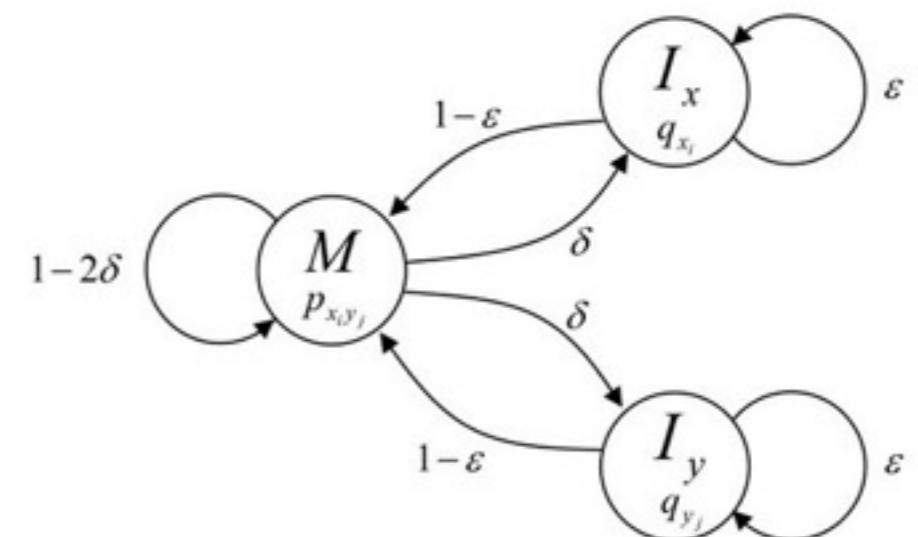
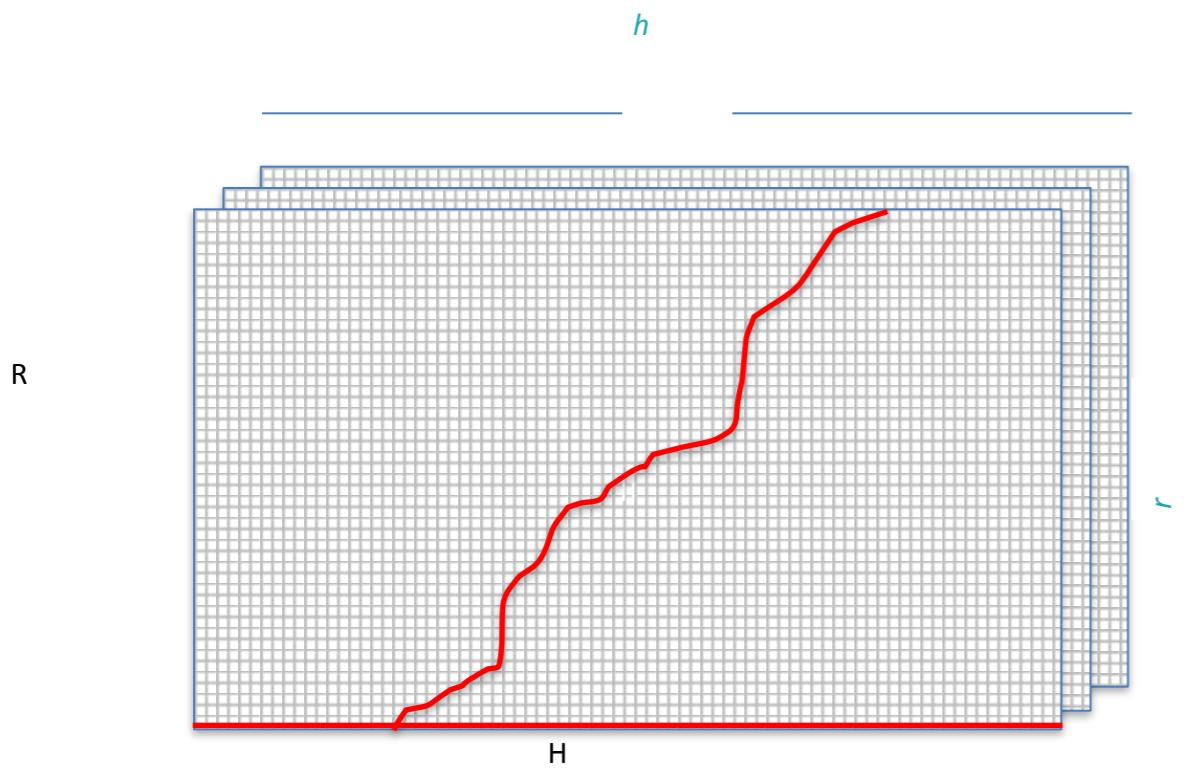


first implementation has shown a 4x speedup but can be improved

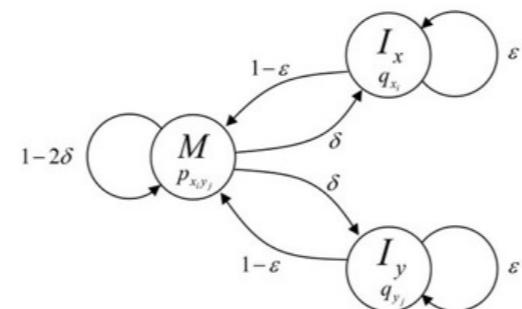
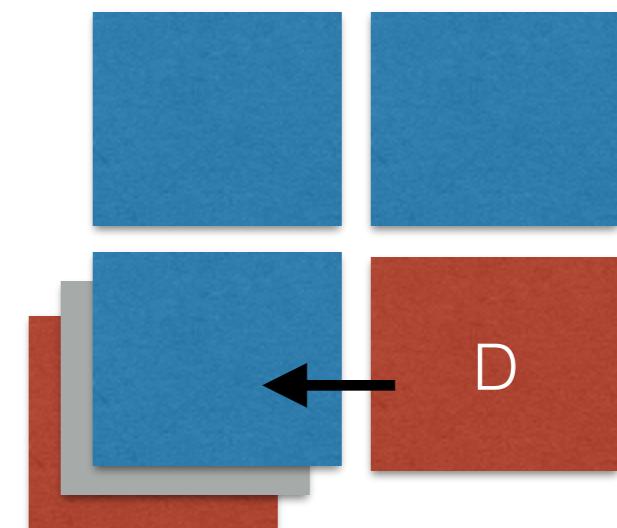
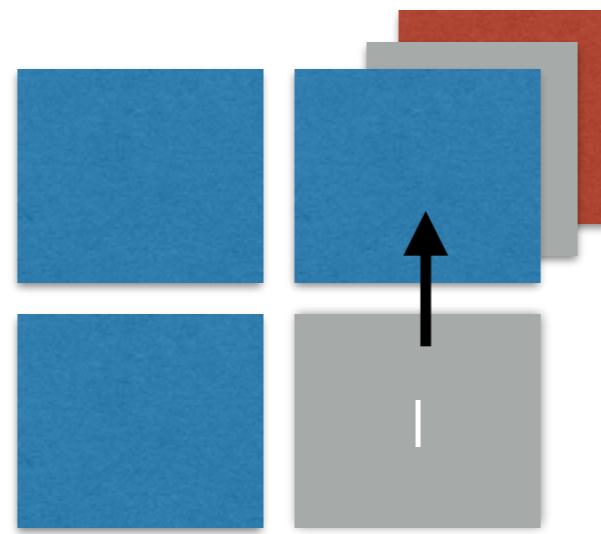
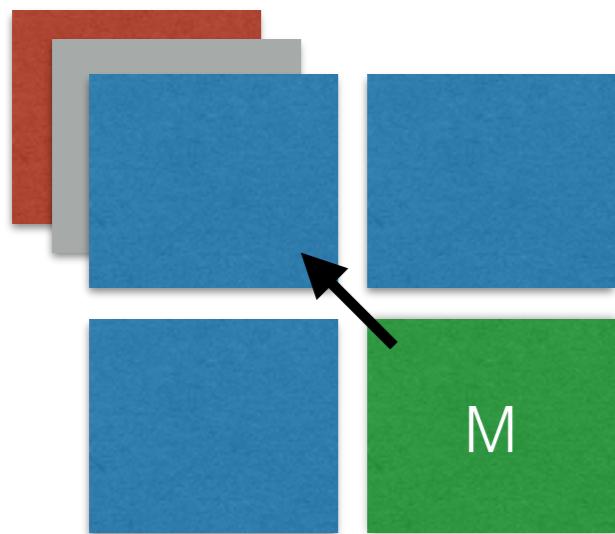
How can we improve the Haplotype Caller performance

- Distributed parallelism: Queue/MapReduce
- Alternative way to calculate likelihoods
- **Heterogeneous parallel compute**

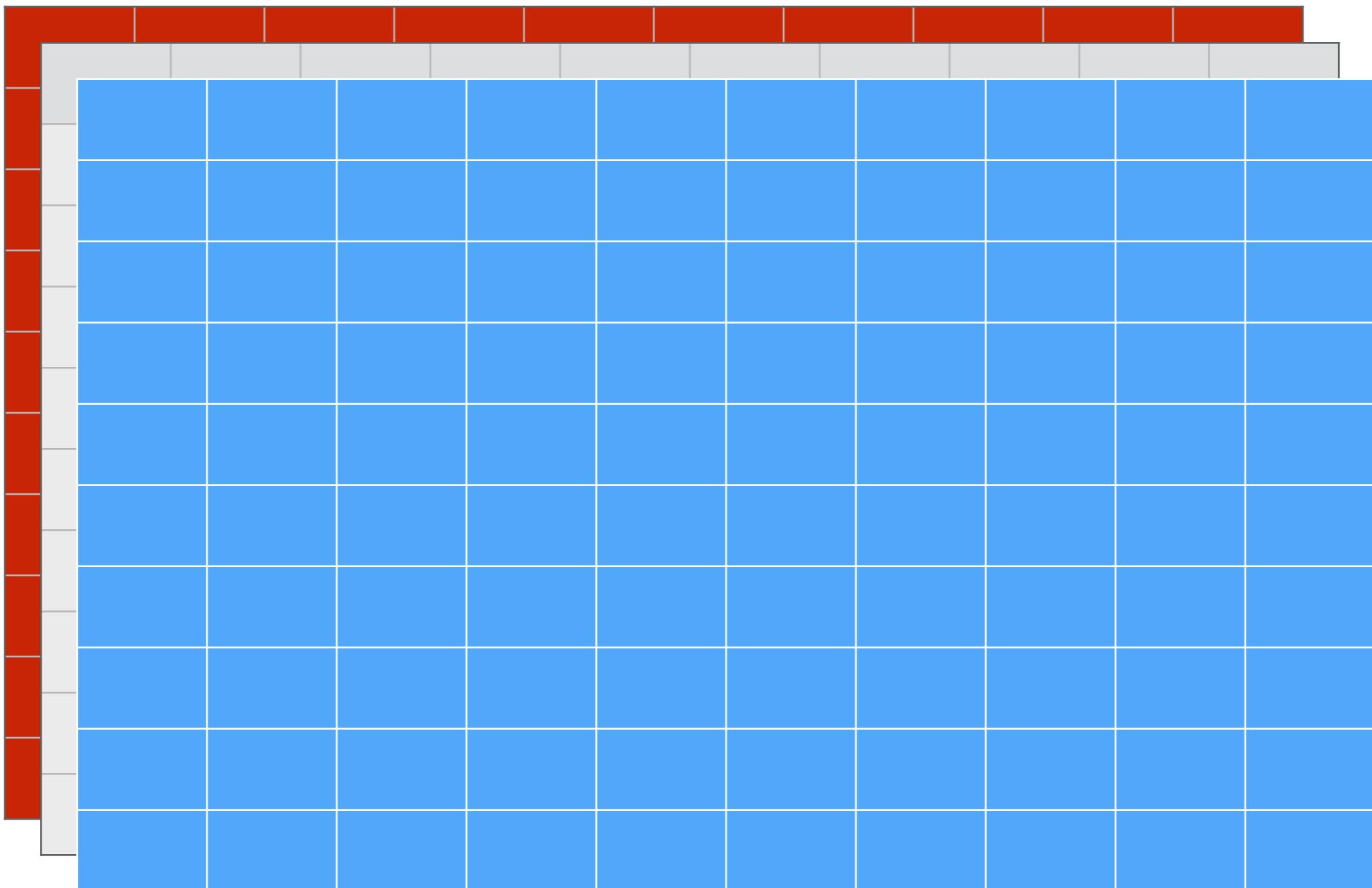
Understanding the Pair-HMM



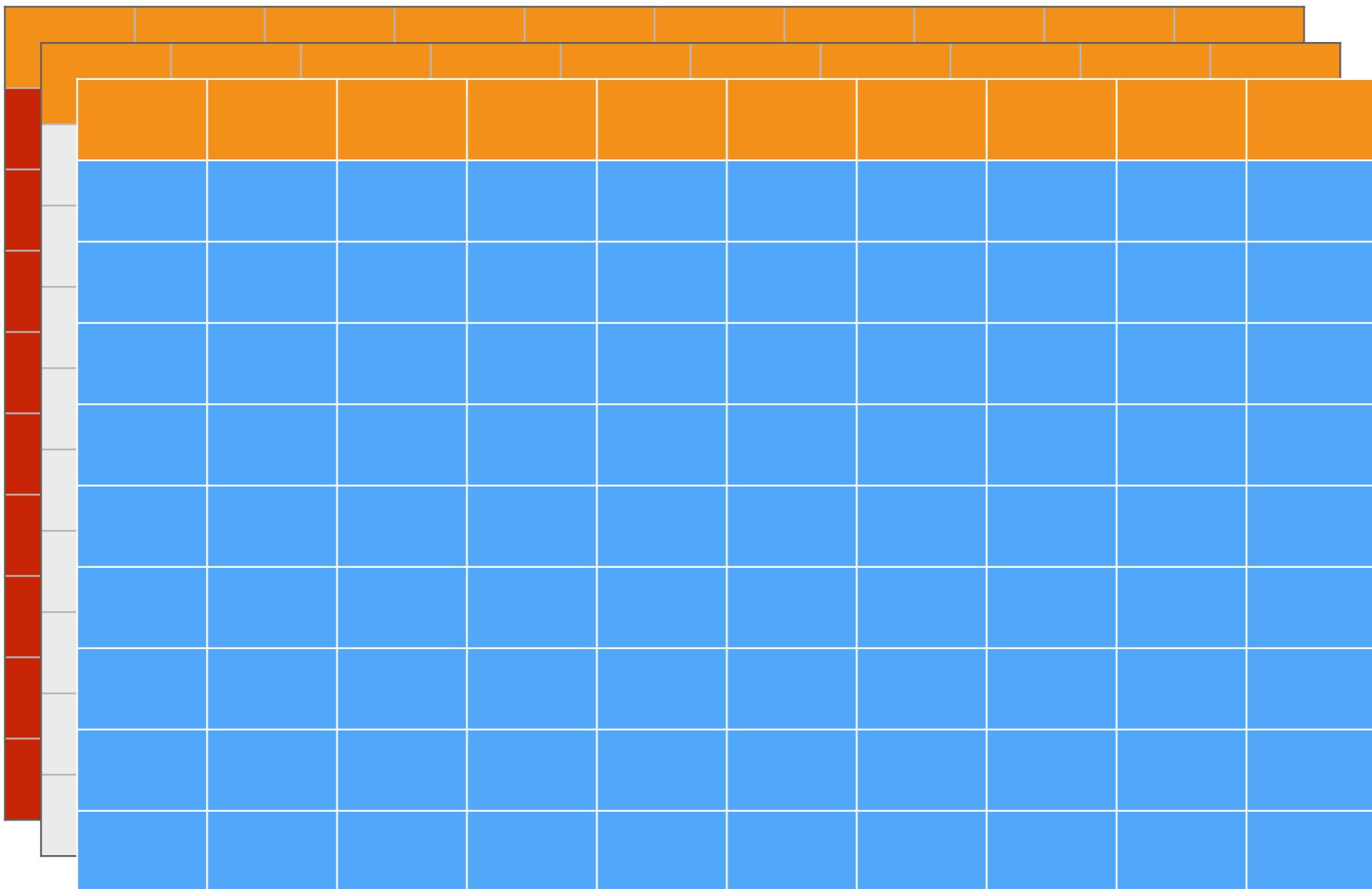
Data dependencies of each cell in each of the three matrices (states)



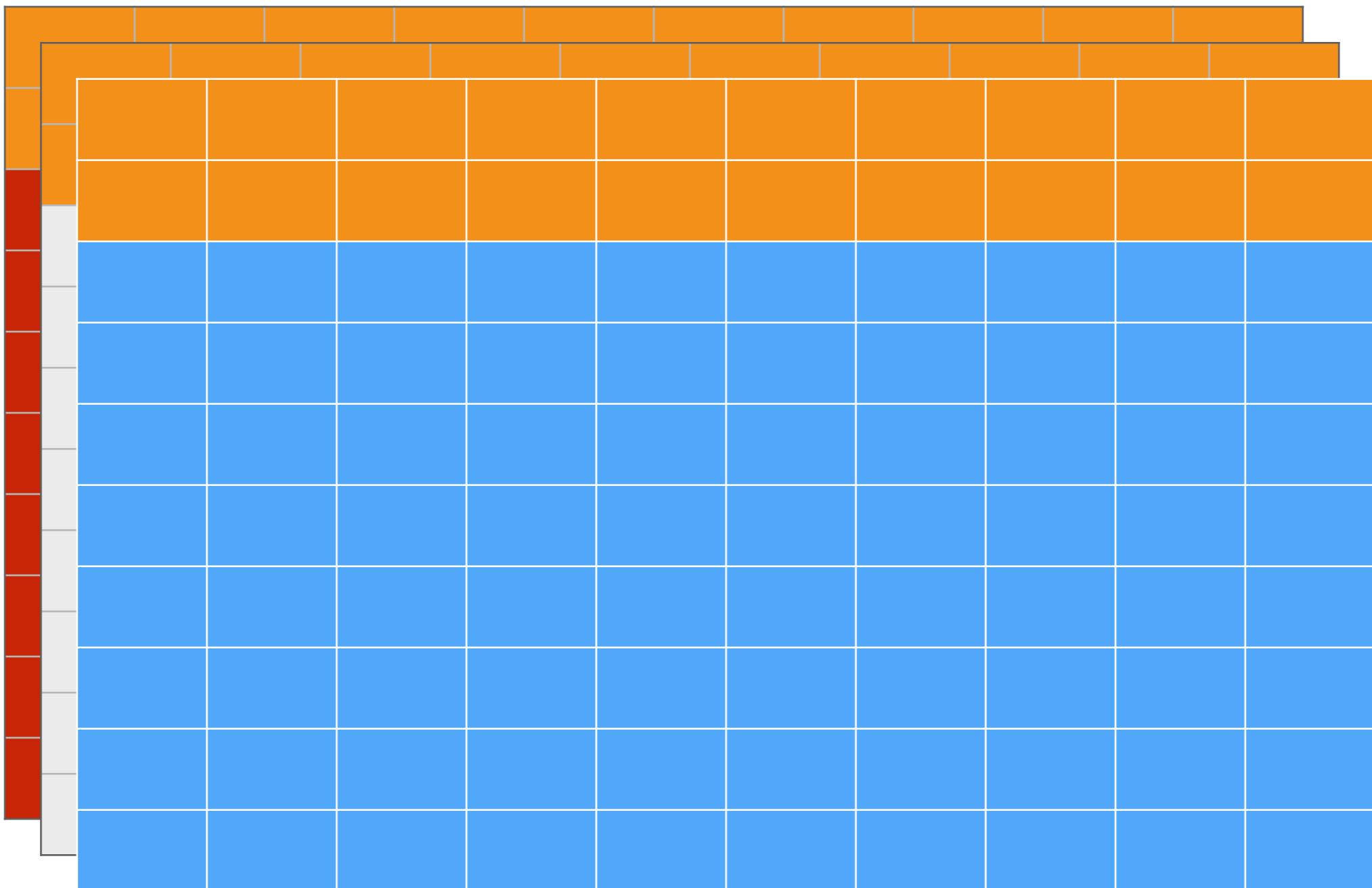
Single core dynamic programming implementation



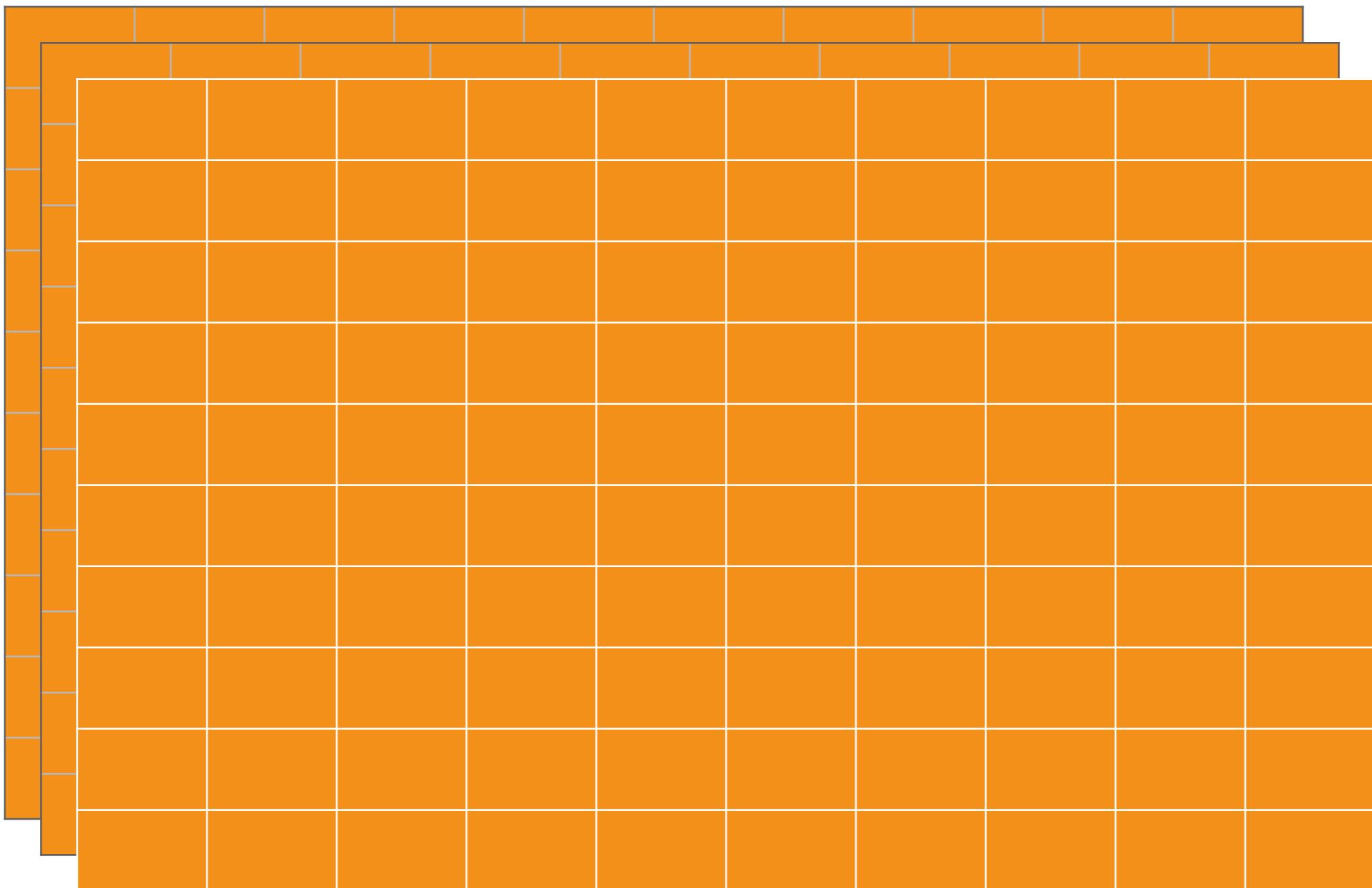
Single core dynamic programming implementation



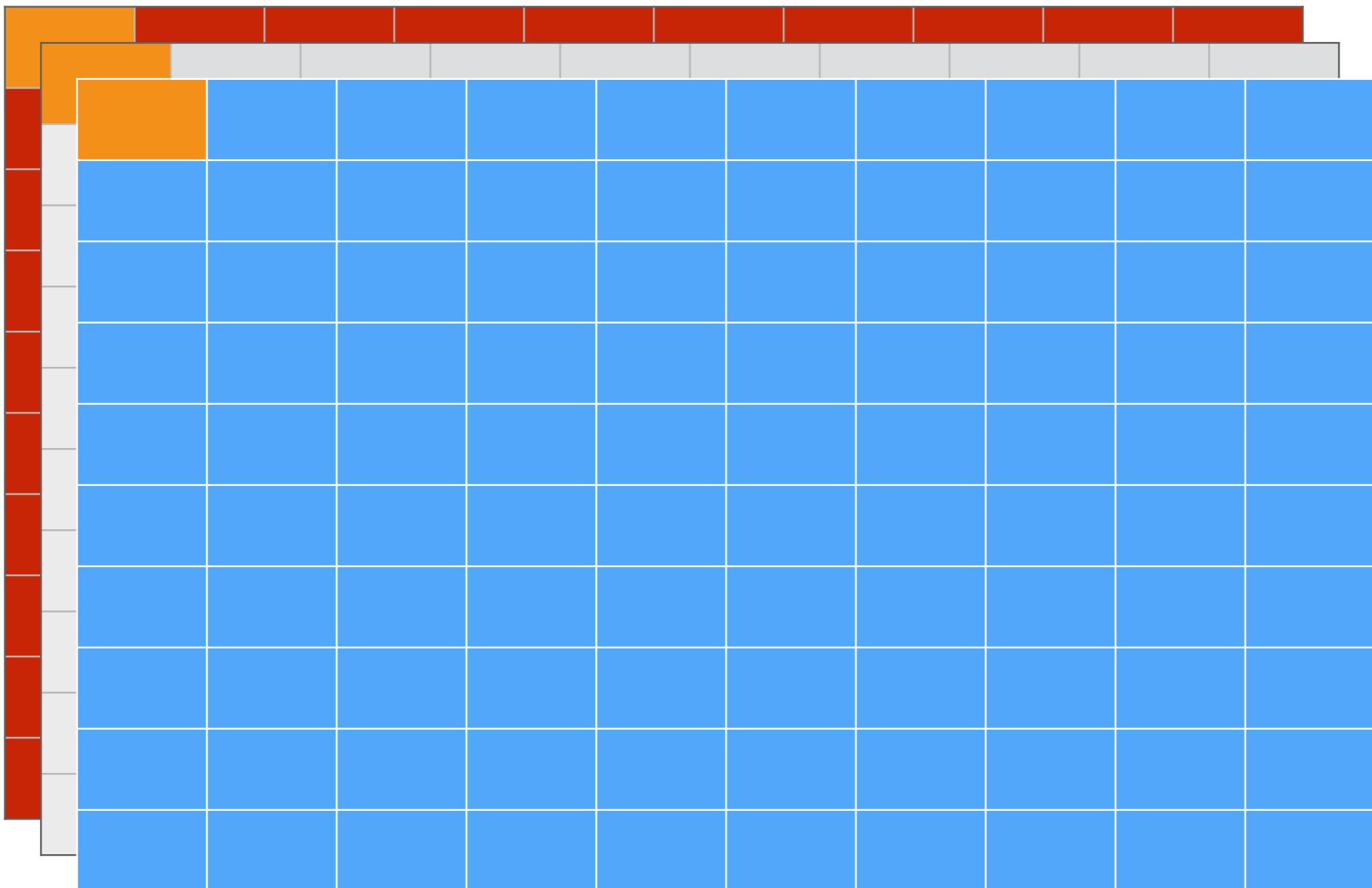
Single core dynamic programming implementation



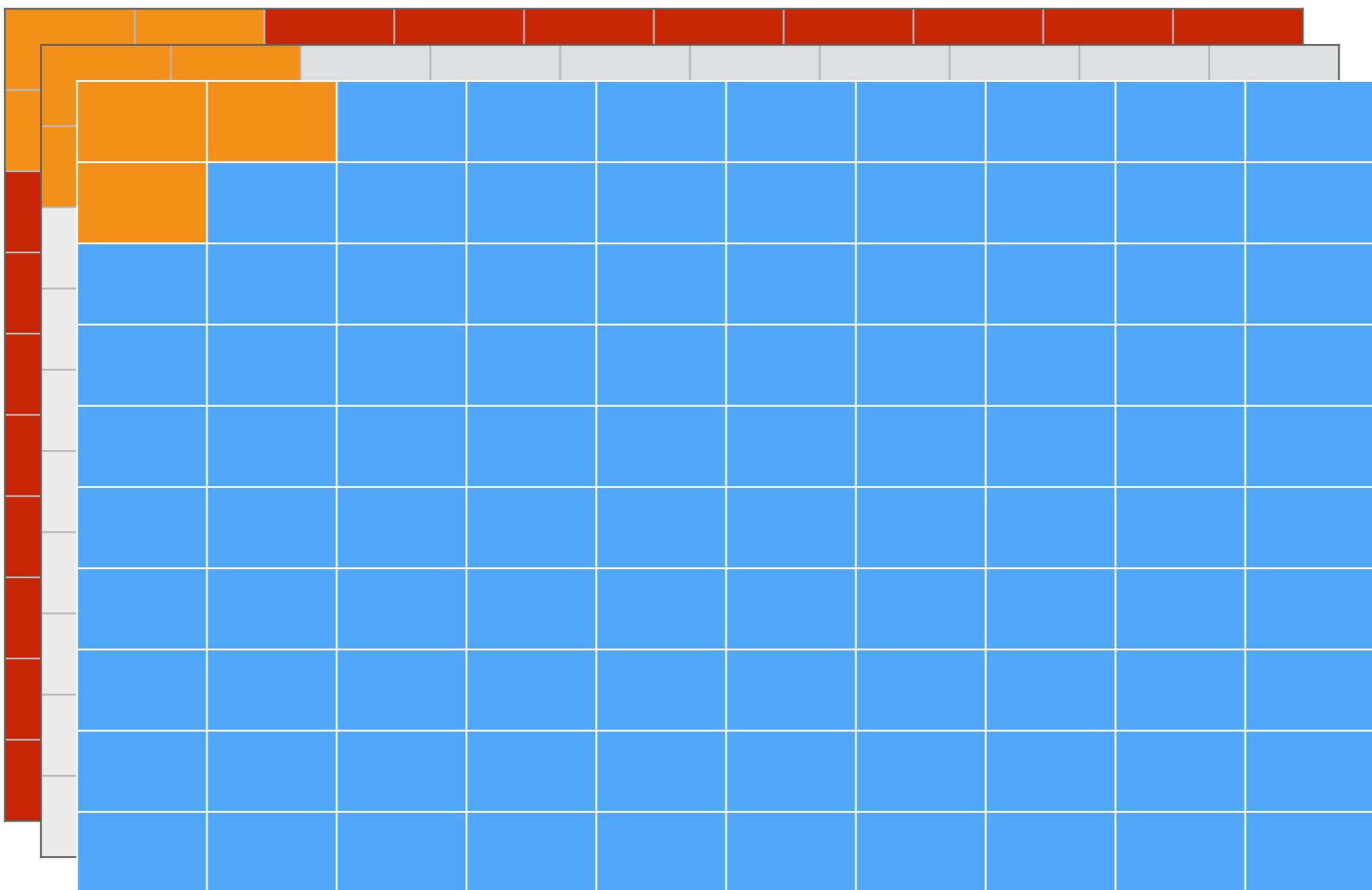
Single core dynamic programming implementation



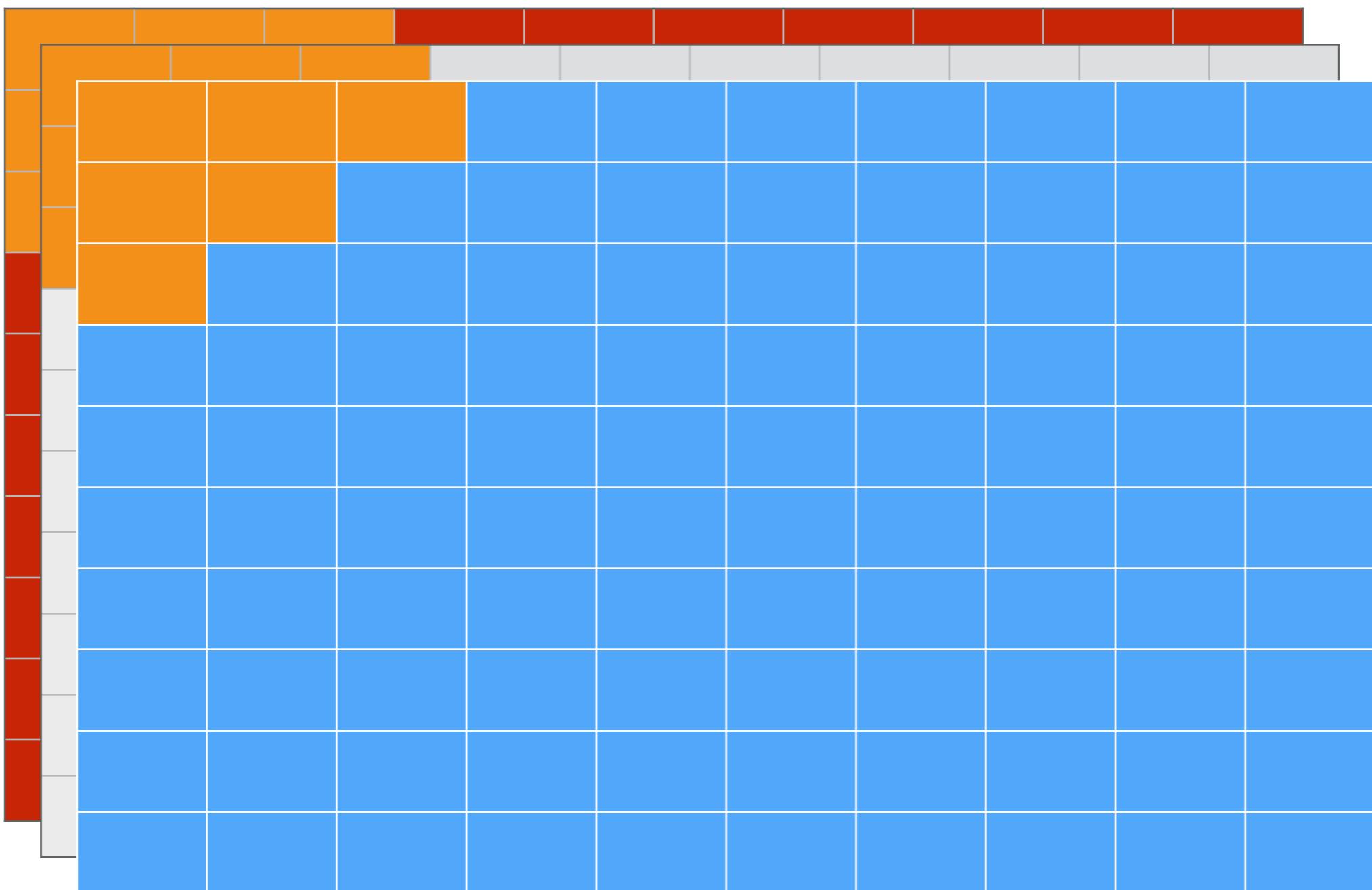
Parallel dynamic programming implementation



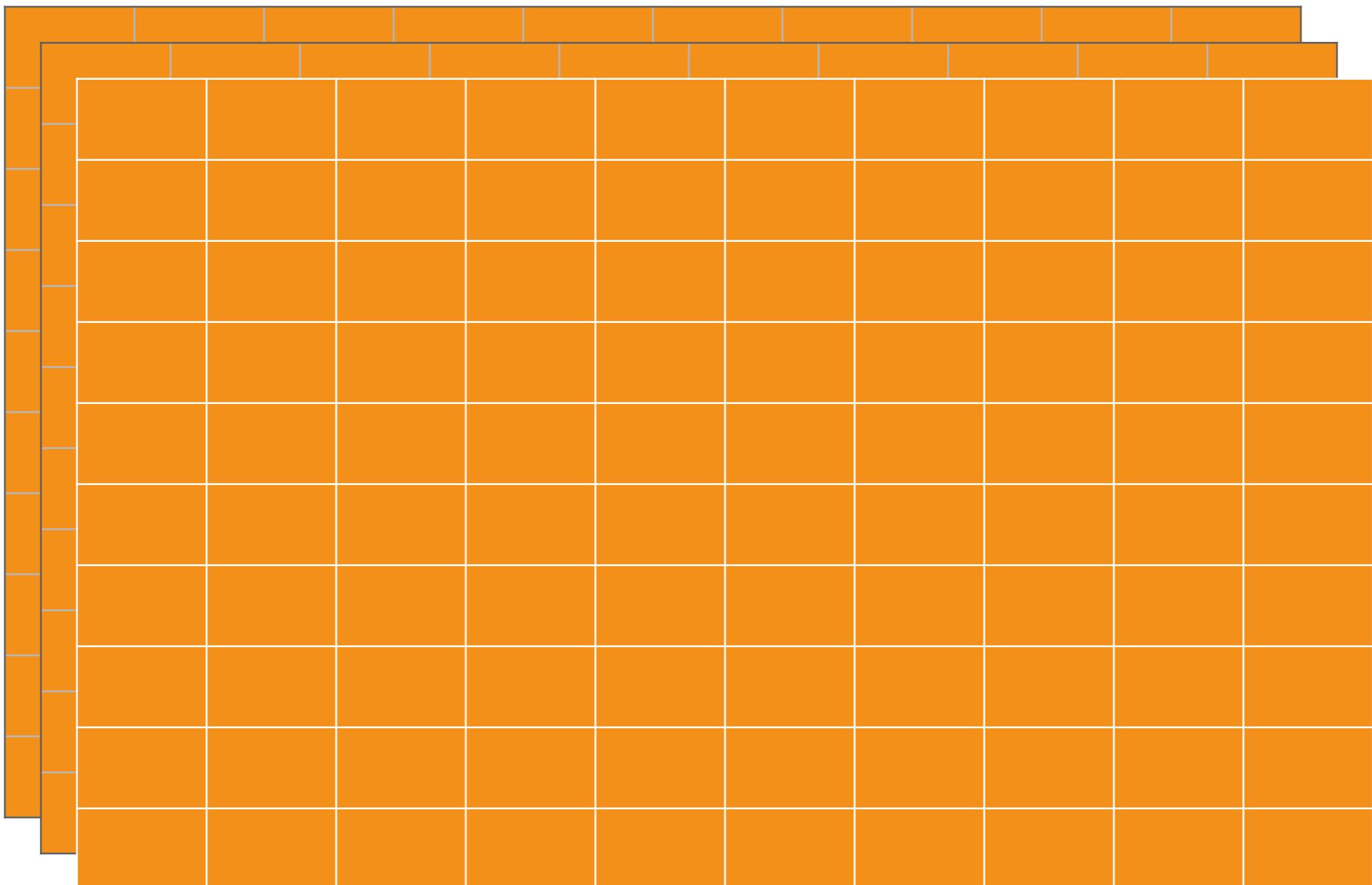
Parallel dynamic programming implementation



Parallel dynamic programming implementation



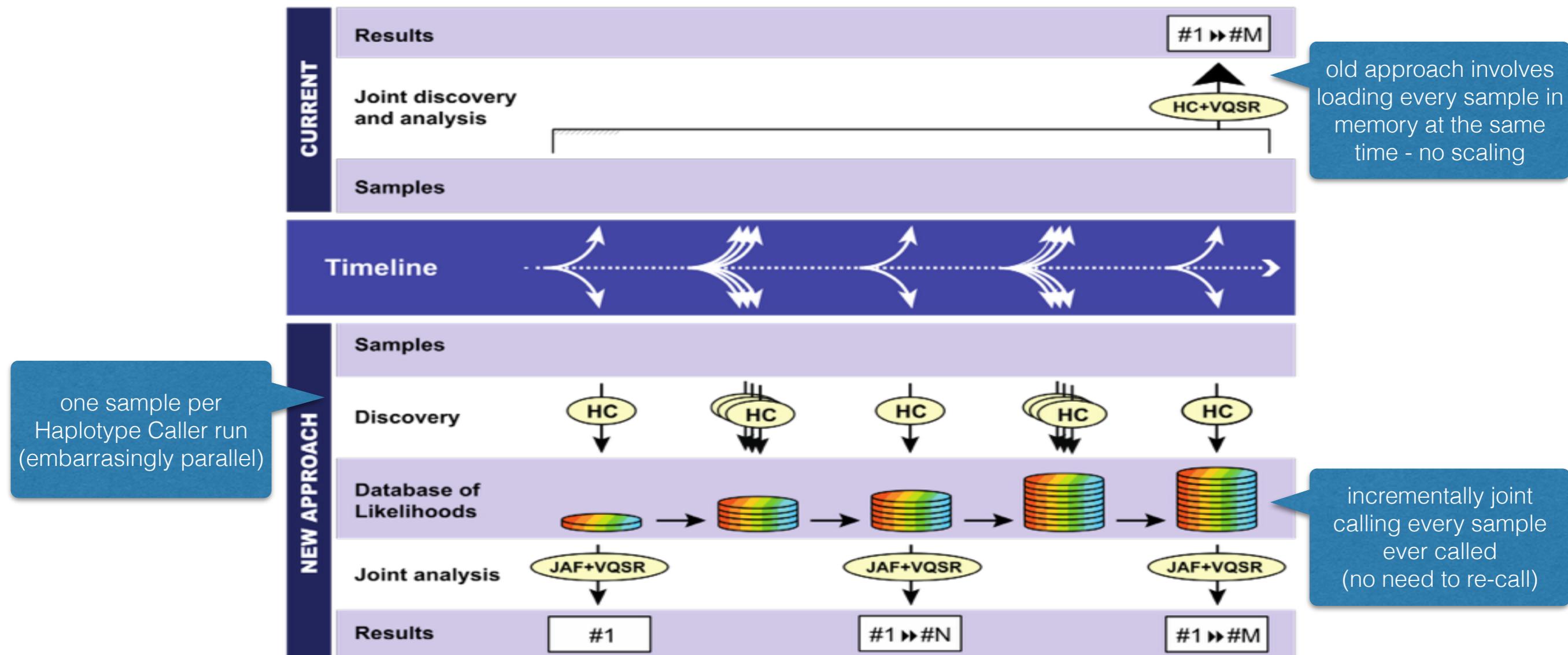
Parallel dynamic programming implementation



Heterogeneous compute speeds up variant calling significantly

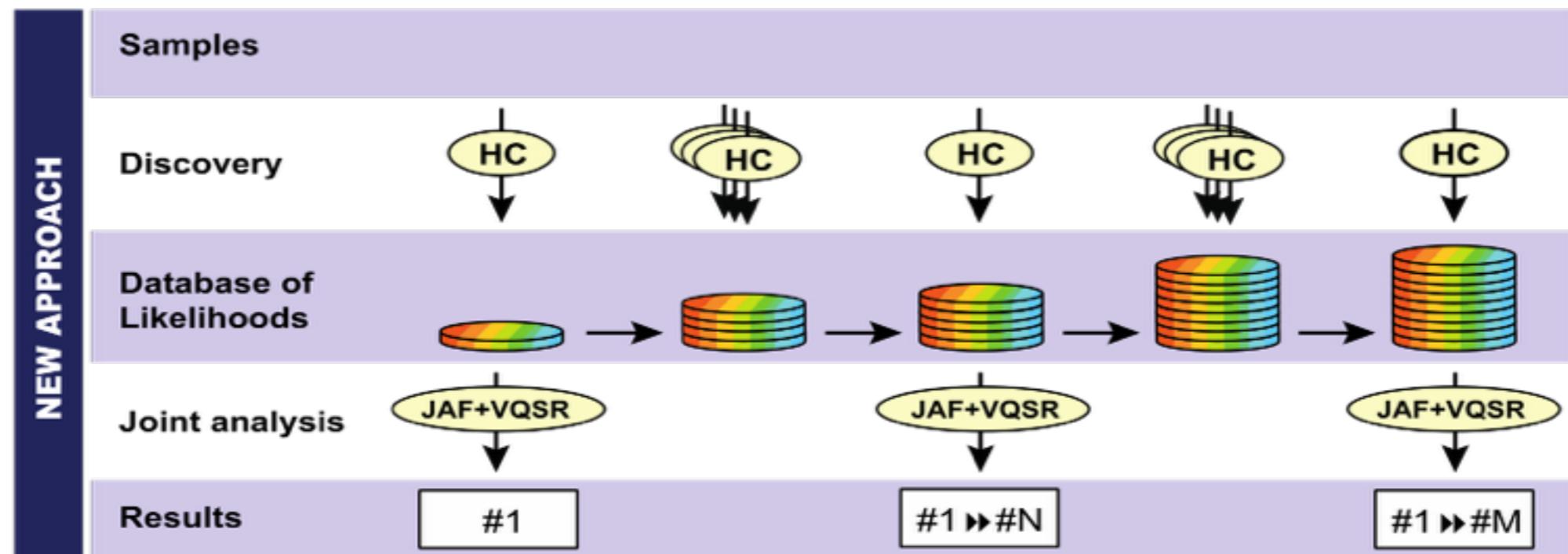
Technology	Hardware	Runtime	Improvement
GPU	NVidia Tesla K40	70	154x
GPU	NVidia GeForce GTX Titan	80	135x
GPU	NVidia GeForce GTX 480	190	56x
GPU	NVidia GeForce GTX 680	274	40x
GPU	NVidia GeForce GTX 670	288	38x
AVX	Intel Xeon 1-core	309	35x
FPGA	Convey Computers HC2	834	13x
-	C++ (baseline)	1,267	9x
-	Java (gatk 2.8)	10,800	-

The reference model enables incremental calling



by separating discovery from joint analysis, we can now jointly call any arbitrary number of samples

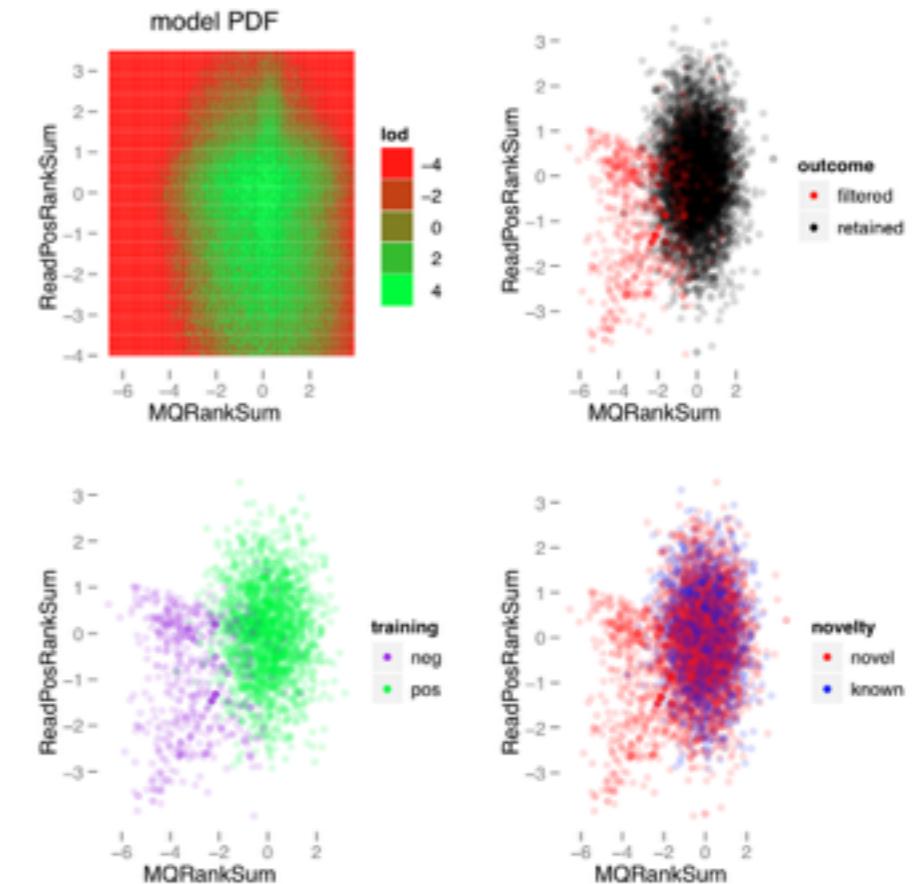
But the joint analysis needs better infrastructure to become a trivial step



currently joint analysis has to be done hierarchically and takes approximately the same amount of time as the discovery step

Future challenges to scale up the variant calling pipeline

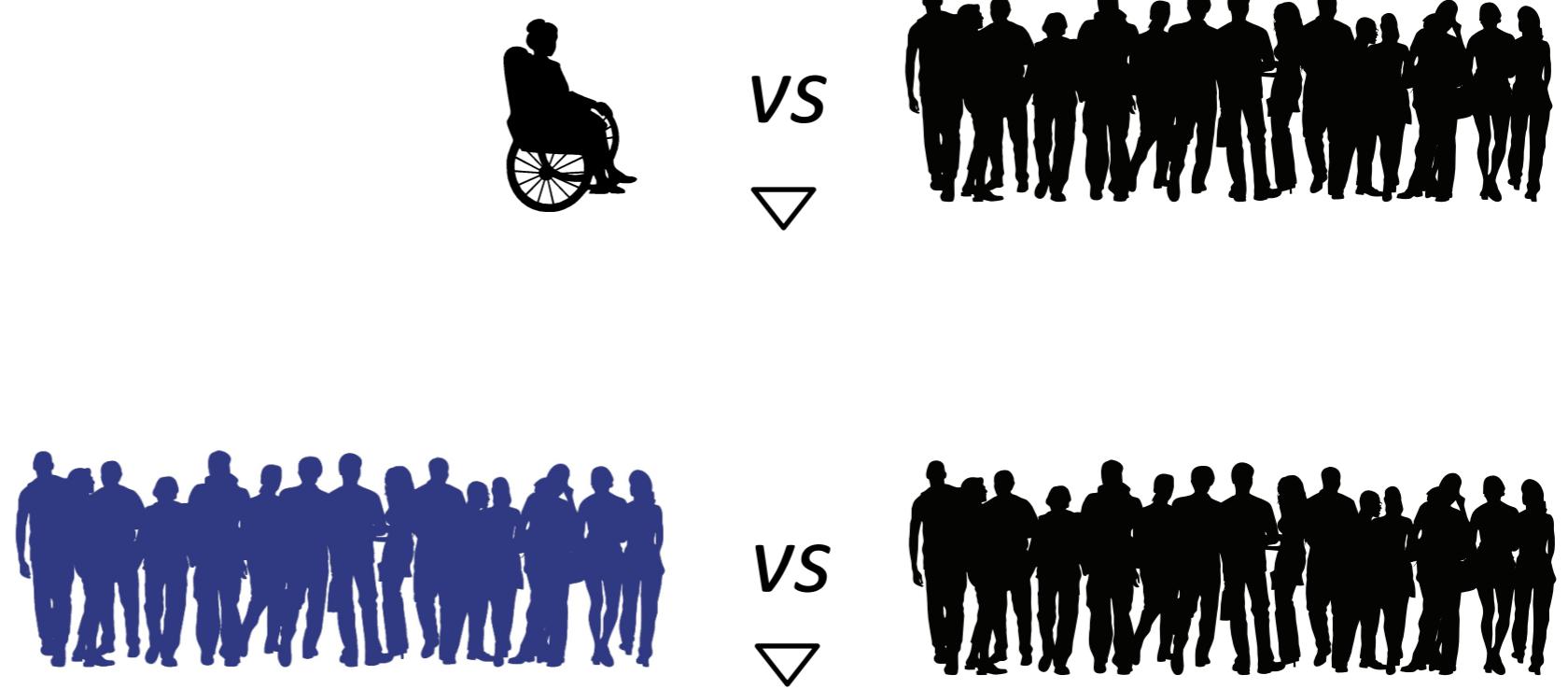
- Improve runtime of the Haplotype Caller (e.g. faster assembly, better prediction of haplotypes)
- Expand on our experience with heterogeneous compute. (e.g. specialize sections of the code for GPU/AVX)
- Improve the performance of the joint analysis by building dedicated high performance tools.
- Improve statistical filtering of indels using non gaussian based methods. (e.g. random forests or other machine learning techniques)
- Revisit all statistical annotations (they behave differently at this scale) identify differences and implement new alternatives.



To fully understand **one** genome we need **tens of thousands** of genomes

Rare Variant
Association Study
(RVAS)

Common Variant
Association Study
(CVAS)



Technical challenge
data integration

Massive data aggregation: The future of large-scale medical sequencing

Proposition

- cost of sequencing has fallen one-million-fold, enabling an explosion of information about the genetic basis of disease
- Learning from the world's combined genomic and clinical data will dramatically accelerate progress
- Aggregated sequence data will be needed to guide the interpretation of genome sequences in clinical practice

Technical challenge

- Aggregating hundreds of thousands of variants and patients requires a sophisticated database system with proper data protection, fast distributed access and a standardized API for tool development.
- For a truly global reach new protective legislation on consents, data access and sharing, IRB processes and patient education will be necessary.

Post-processing analysis pipeline standardization and scaling

- What happens after our current pipeline is extremely undefined.
 - Hundreds of completely unrelated tools are chained together with non-reusable scripts.
 - Analyses are very often non-repeatable.
 - Tools are not generalized and performance does not scale. (usually written in matlab, R, PERL...)
 - Most tools are written by one grad student/postdoc and is no longer maintained.
 - Complementary data types are not standardized (e.g. phenotypic data).

gmap gmod mrfast circos soap
amos pacbio statgen snpeff PAGIT pindel
clcbio genomeview htseq VarScan shrimp picard
tabix bedtools Bioconductor dnnaa
vague VCFtools rnaseq Megan circoletto cutadapt
FastQC solexaqa RepeatMasker bamtools allpaths-lg genepattern
glimmer VariationHunter smash community Mapper rdp_classifier
SnipSift novelseq gatk metagene smalt bwa smrtanalysisrstudio
mauve breakaway GAPipeline browse seqclean varsifter sra SVMerge
tophat breakdancer assembler bwa rdxplorer velvet
glimmer prinseq jbrowse seqclean SFFtools
bowtie celera cctools.snpEff Assembler pasha SVDetect Tablet
taverna Quime gsAssembler samtools Savant
gmapexonerate metasim Hydra Mira dindel
cnv-seq

Future challenges to enable large scale post-processing analysis

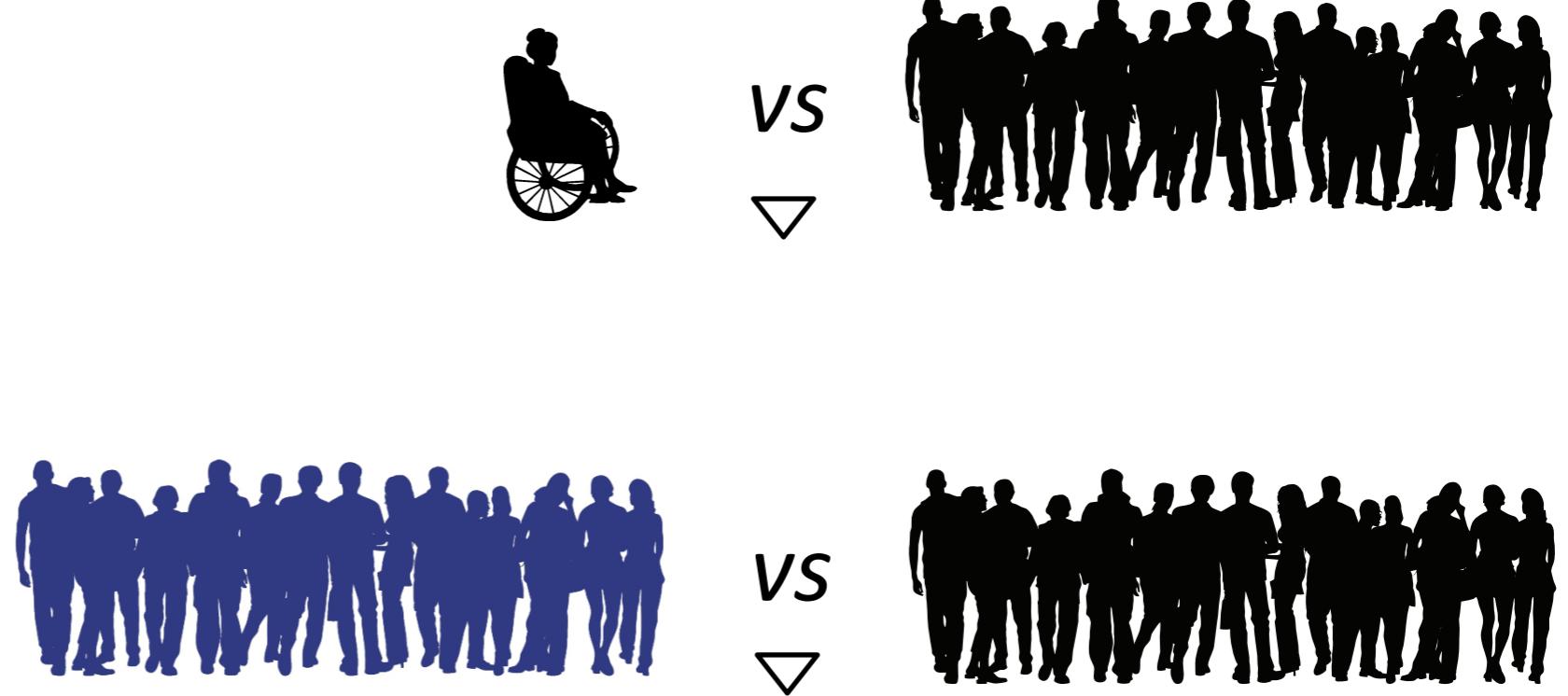
- Build a generic *variant database* that would enable large scale data analysis of many thousands of patients and controls
- Establish best practices and revisit the available tools to create a uniform framework where developers can create reusable and maintainable tools.
- Standardize workflows for common analyses such as RVAS and CVAS.
- Establish best practices for use of such pipelines in the clinic.



To fully understand **one** genome we need **tens of thousands** of genomes

Rare Variant
Association Study
(RVAS)

Common Variant
Association Study
(CVAS)



Technical challenges
RNA-seq, Structural Variation and Cancer

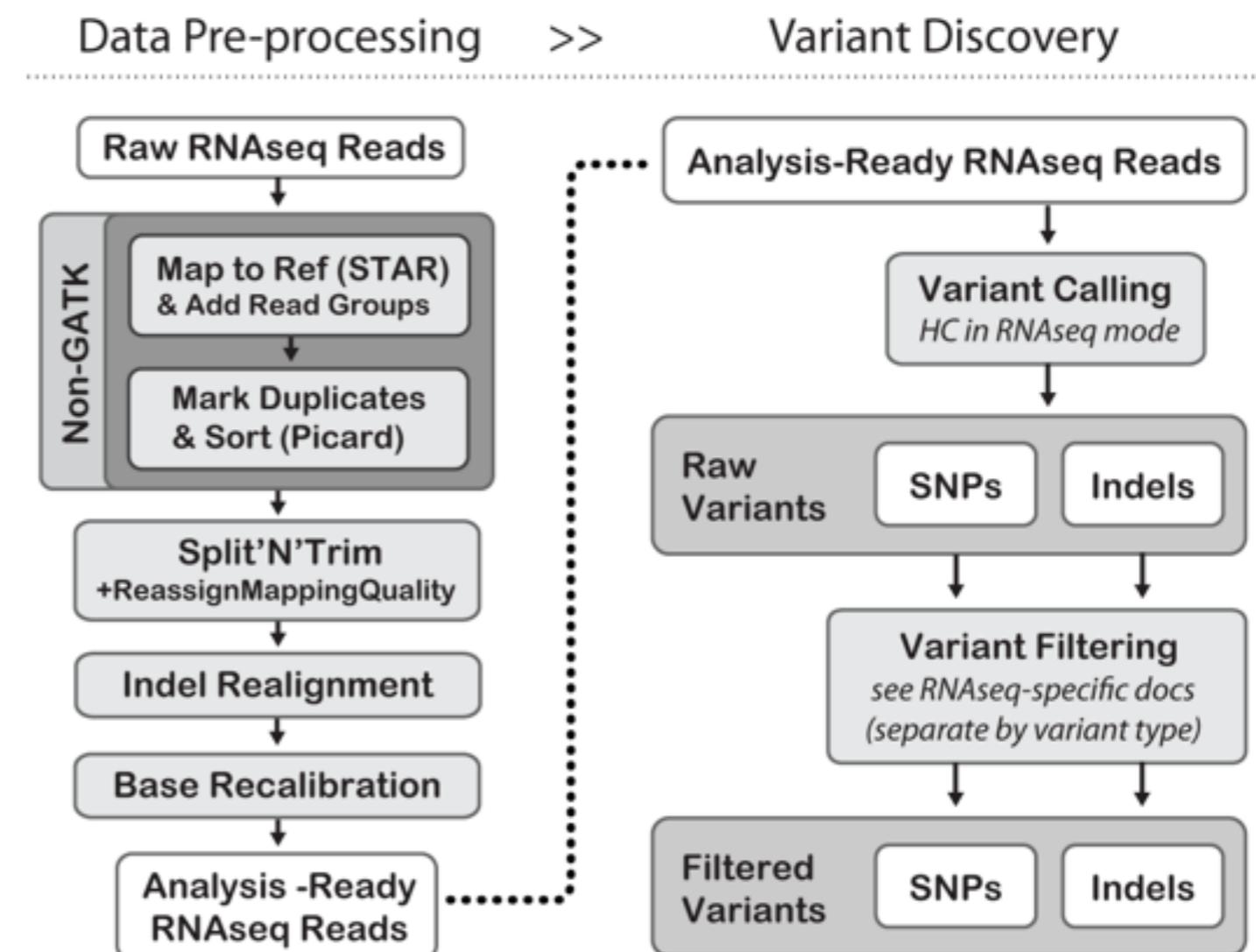
DNA does not tell the whole story — there is RNA too!

Milestone 1:

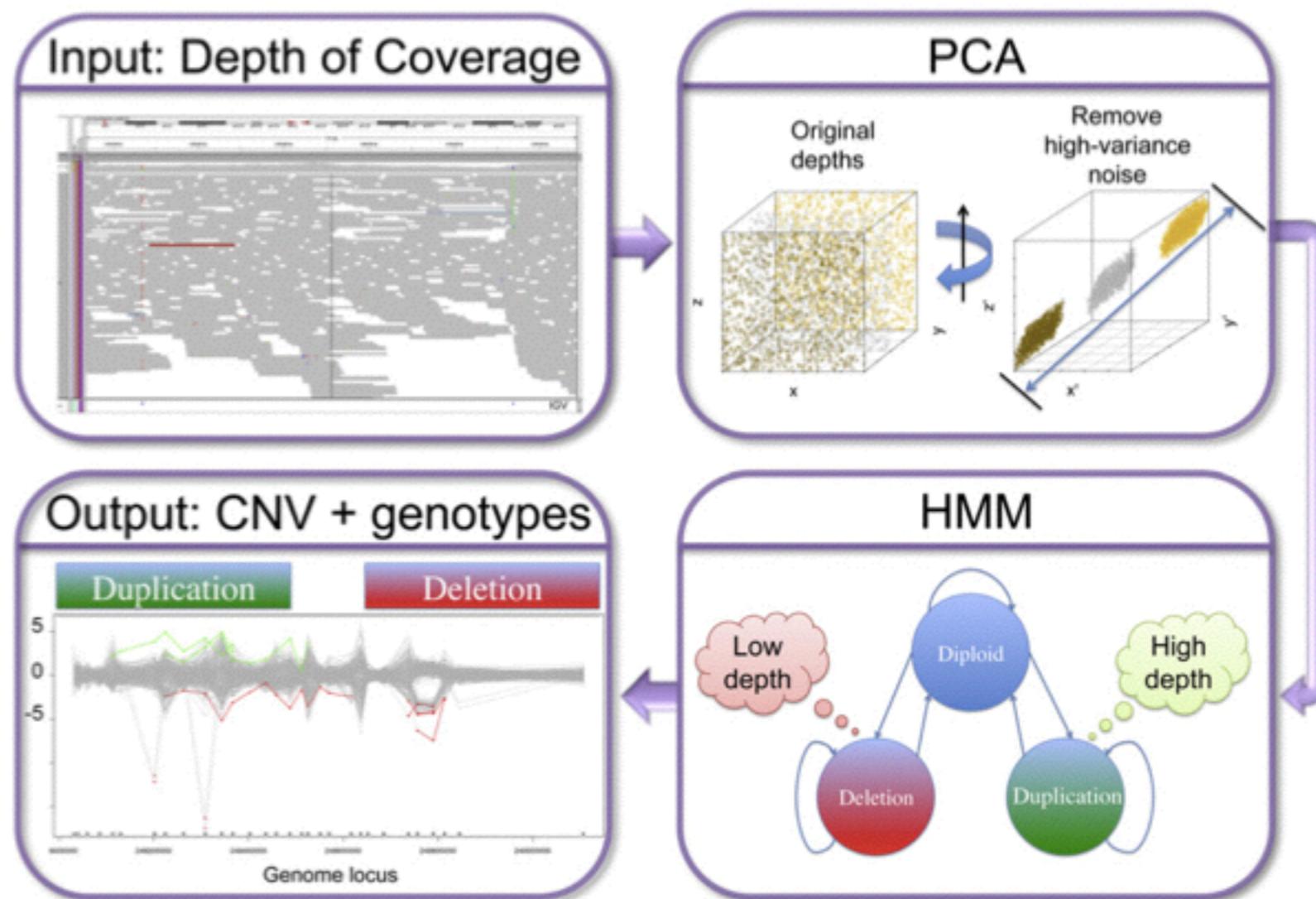
update GATK tools to make best use of RNA data including contrastive variant calling with DNA. Improve accuracy and overall performance.

Milestone 2:

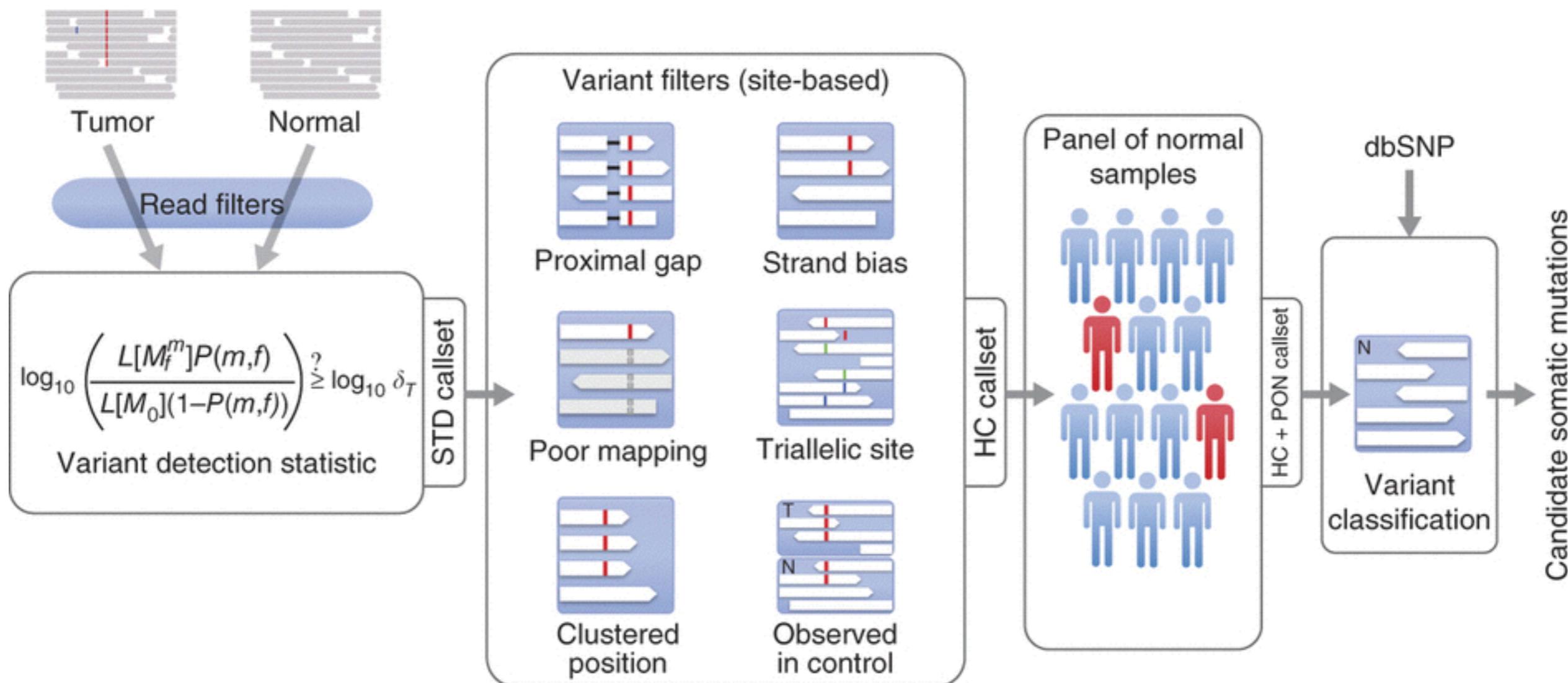
Build new tools to address the needs of the RNA-seq world.
(e.g. Heterogeneous expression, RNA editing, expression levels, alternate splicing...)



Structural variation is an important missing piece in the analysis pipeline

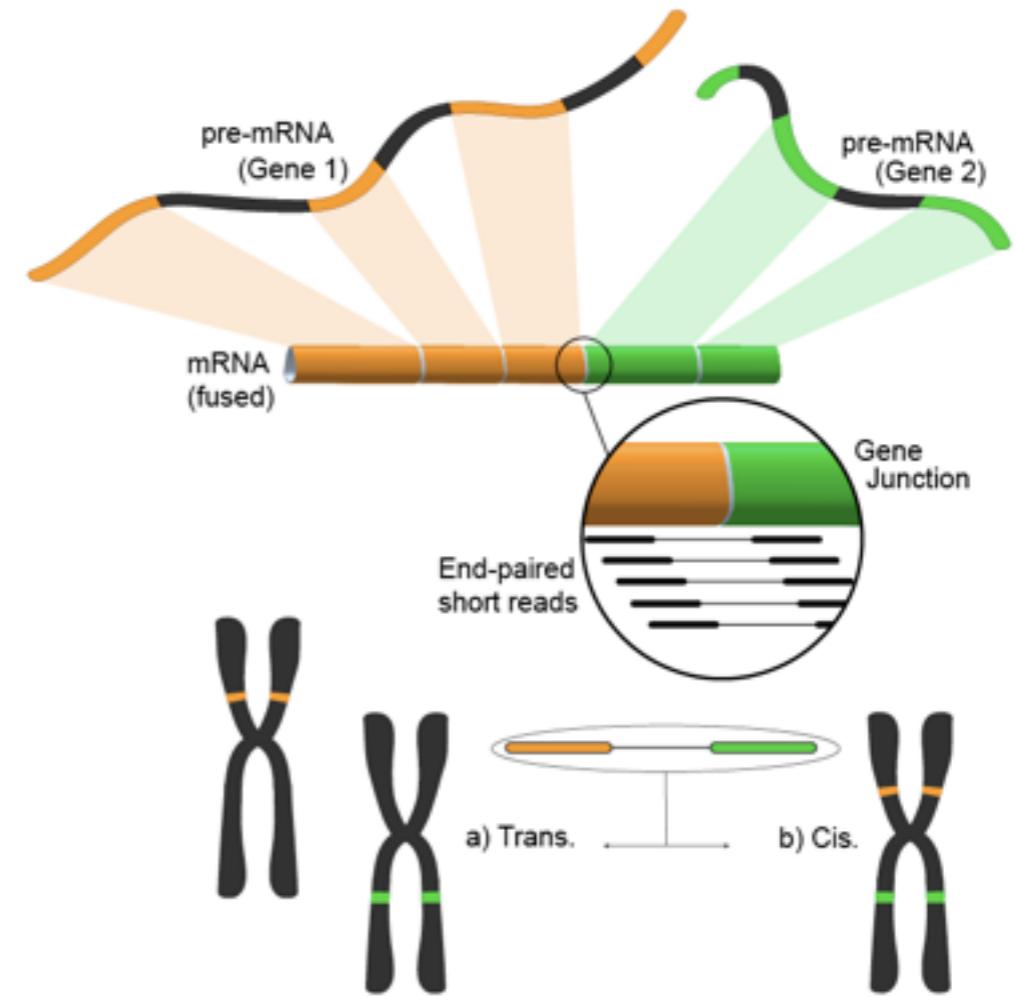


Cancer needs the same rigorous standardization from study design to data processing and analysis



Future challenges to scale up the variant calling pipeline

- RNA-seq has great potential to further complement discoveries from DNA-seq.
- The world of clinical applications has yet to explore the plethora of information provided by the RNA — we need tools!
- Structural variations are known to play an important role. Tools should be integrated into the standard research pipeline for CVAS and RVAS.
- Cancer still has many open questions from the processing side to the specialized downstream analysis.



There is a lot to do... but it's definitely worth it.

- samples must be consistently pre-processed and the processing pipelines need to scale in performance.
- Variants must be jointly called and currently available tools need to provide the necessary performance. (We solved the scaling problem!)
- Data needs to be integrated to vastly accelerate discovery. New technology needs to be built in order to enable the creation of methods and tools that leverage the aggregated data.
- RNA-seq and other analyses need to be integrated and standardized for scientists and clinicians to understand the whole picture.

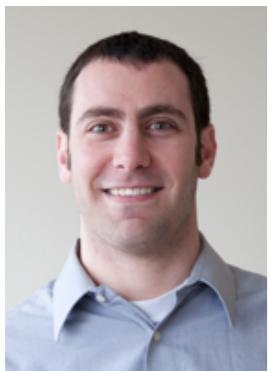
Rare Variant
Association Study
(RVAS)

Common Variant
Association Study
(CVAS)



This is the work of many...

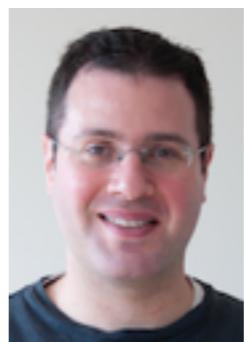
team



colleagues



collaborators





“Thank you!”

–Mauricio Carneiro