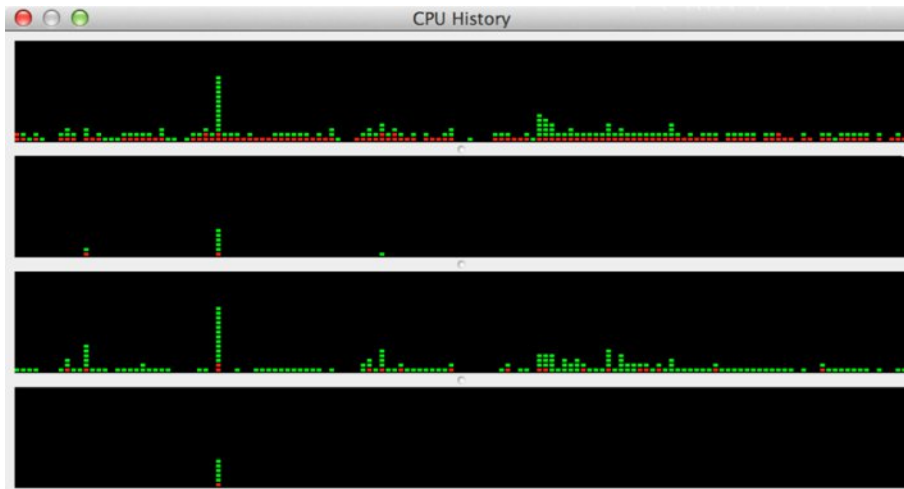
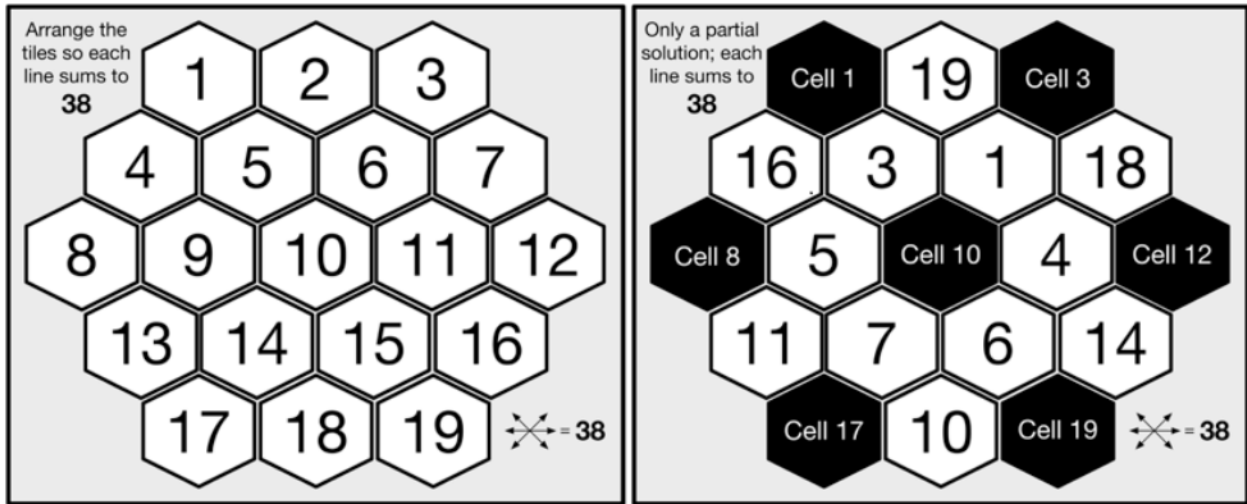


Chapter 1: Simple Parallelism with R



Activity Monitor (All Processes)

CPU Memory Energy Disk Network

Process Name	% CPU	CPU Time	Threads	Idle Wake Ups	PID	User
rsession	95.6	26:37.97	1	0	38883	simon
rsession	95.6	26:38.67	1	0	38881	simon
rsession	95.5	26:38.02	1	0	38884	simon
rsession	94.2	26:38.29	1	0	38882	simon
Activity Monitor	4.1	46:30.57	4	3	7694	simon

Activity Monitor (My Processes)

CPU Memory Energy Disk Network

Process Name	% CPU	CPU Time	Thr...	Idle Wake Ups	PID	User
R	0.0	2.25	1	0	41527	simon
R	0.0	1:22.98	1	0	41551	simon
R	0.0	1:22.33	1	0	41543	simon
R	0.0	55.94	1	0	41535	simon

Process Name	% CPU	CPU Time	Thr... ▲	Idle Wake Ups	PID	User
R	0.0	6:34.15	1	0	42720	simon
R	0.0	3:34.99	1	0	42712	simon
R	0.0	8.36	1	0	42704	simon
R	0.0	7:26.65	1	0	42728	simon

Process Name	% CPU	CPU Time	Thr...	Idle Wake Ups	PID	User
R	0.0	6:51.47	1	0	43092	simon
R	0.0	6:14.22	1	0	43084	simon
R	0.0	6:12.69	1	0	43076	simon
R	0.0	5:14.08	1	0	43100	simon

Create Access Key ✕

✔ Your access key (access key ID and secret access key) has been created successfully.

Download your key file now, which contains your new access key ID and secret access key. If you do not download the key file now, you will not be able to retrieve your secret access key again.

To help protect your security, store your secret access key securely and do not share it.

[▶ Show Access Key](#)

Download Key File
Close

Services ▼ Edit ▼
Simon ▼ N. Virginia ▼ Help ▼

Elastic MapReduce ▼ Cluster List
EMR Help

Create cluster
View details
Clone
Terminate

Filter: All clusters ⌵ Filter clusters ...
13 clusters (all loaded)

	Name	ID	Status	Creation time (UTC+1) ▼	Elapsed time	Normalized instance hours
▶ ●	RJob-Sun Oct 19 20:50:11 2014	j-2JTTPF0MAPYHF	Waiting	2014-10-19 20:50 (UTC+1)	41 minutes	64

Summary

Master ec2-54-164-238-82.compute-public DNS: 1.amazonaws.com

Termination protection: Off Change

Tags: -- View All / Edit

Hardware Resize

Master: Running 1 m1.large

Core: Running 15 m1.large

Task: --

[View cluster details](#)

Add Step			
Name	Status	Start time (UTC+1) ▼	Elapsed time
2014-10-19_20:56:25.16105	Completed	2014-10-19 20:56 (UTC+1)	34 minutes

Bootstrap Actions

Name
R-InstallLatest
R-UpdatePackages

Services Edit Simon N. Virginia Help

EC2 Dashboard

Launch Instance Connect Actions

Filter by tags and attributes or search by keyword 1 to 36 of 36

Instance ID	Instance Type	Availability Zone	Instance State	Status Checks	Alarm Status	Public DNS	Pub
i-2d8466c7	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-170-64.co...	54
i-3a8466d0	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-176-213.co...	54
i-388466d2	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-181-141.co...	54
i-398466d3	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-180-23.co...	54
i-3f8466d5	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-171-205.co...	54
i-3c8466d6	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-181-41.co...	54
i-3d8466d7	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-165-246.co...	54
i-328466d8	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-187-199.co...	54
i-338466d9	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-187-230.co...	54
i-308466da	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-188-219.co...	54
i-318466db	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-174-139.co...	54
i-308466dc	m1.large	us-east-1c	running	2/2 checks ...	None	ec2-54-172-180-123.co...	54

Select an instance above

© 2008 - 2014, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use Feedback

https://console.aws.amazon.com/console/home?region=us-east-1

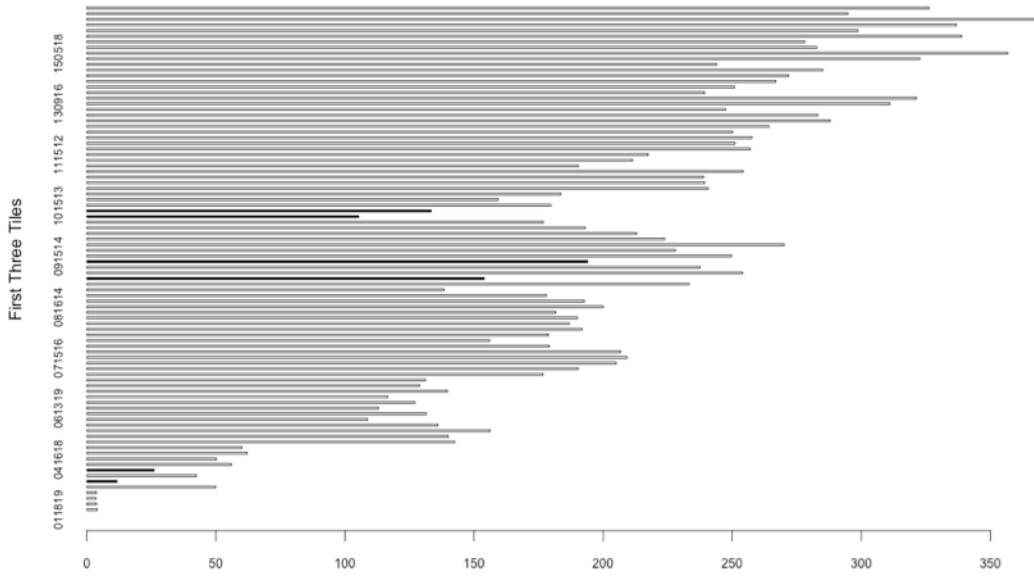
Services Edit Simon Global Help

Upload Create Folder Actions None Properties Transfers

All Buckets / rtmphcnz7ibnqkjlora-segueout / results

Name	Storage Class	Size	Last Modified
part-0000	Standard	2.5 KB	Sun Oct 19 21:07:58 GMT+100 2014
part-0001	Standard	2 KB	Sun Oct 19 21:04:11 GMT+100 2014
part-0003	Standard	3.1 KB	Sun Oct 19 21:12:44 GMT+100 2014
part-0004	Standard	1.4 KB	Sun Oct 19 21:03:25 GMT+100 2014
part-0006	Standard	3.1 KB	Sun Oct 19 21:08:59 GMT+100 2014
part-0007	Standard	1.4 KB	Sun Oct 19 21:04:13 GMT+100 2014
part-0009	Standard	2.6 KB	Sun Oct 19 21:09:04 GMT+100 2014
part-0010	Standard	2 KB	Sun Oct 19 21:06:58 GMT+100 2014
part-0011	Standard	2.5 KB	Sun Oct 19 21:11:45 GMT+100 2014
part-0013	Standard	3.1 KB	Sun Oct 19 21:11:13 GMT+100 2014
part-0014	Standard	2.6 KB	Sun Oct 19 21:13:42 GMT+100 2014
part-0017	Standard	0 bytes	Sun Oct 19 21:04:29 GMT+100 2014
part-0018	Standard	2.6 KB	Sun Oct 19 21:11:07 GMT+100 2014
part-0022	Standard	2.6 KB	Sun Oct 19 21:13:37 GMT+100 2014

AWS EMR Solver Execution Profile



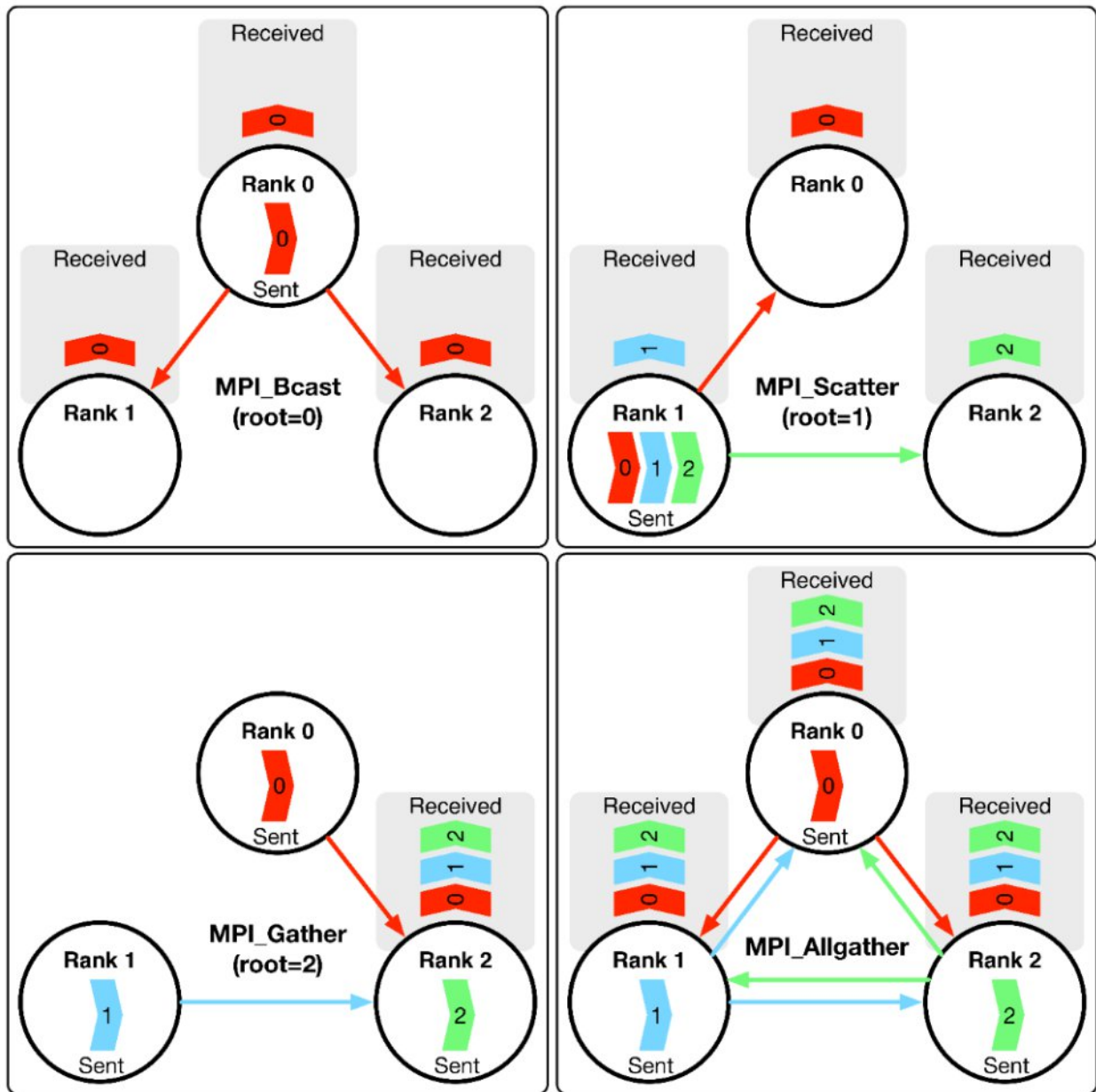
Boards=90 -- Min=0m4s Max=6m8s Avg=3m13s -- Fastest solution 031718 in 12s

Chapter 2: Introduction to Message Passing

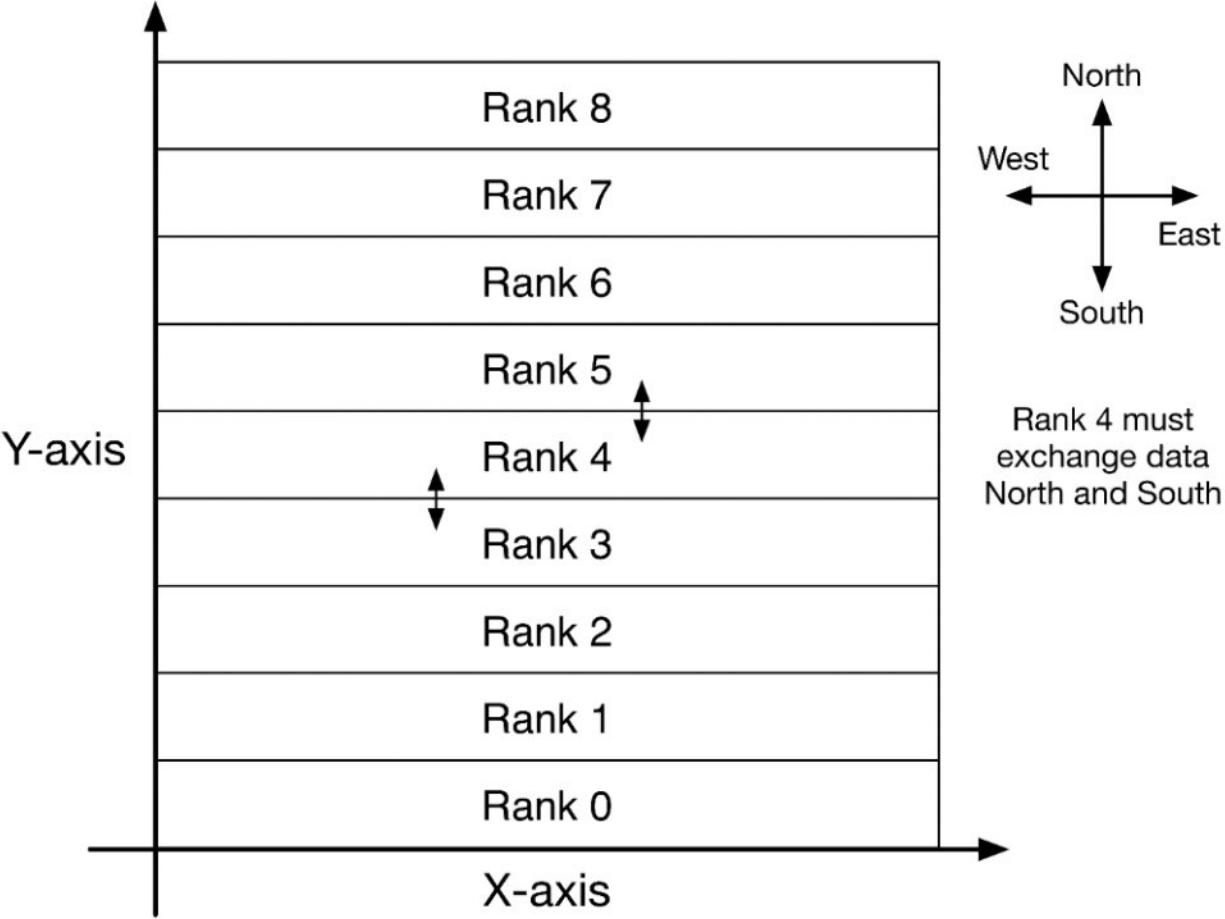
Activity Monitor (All Processes)

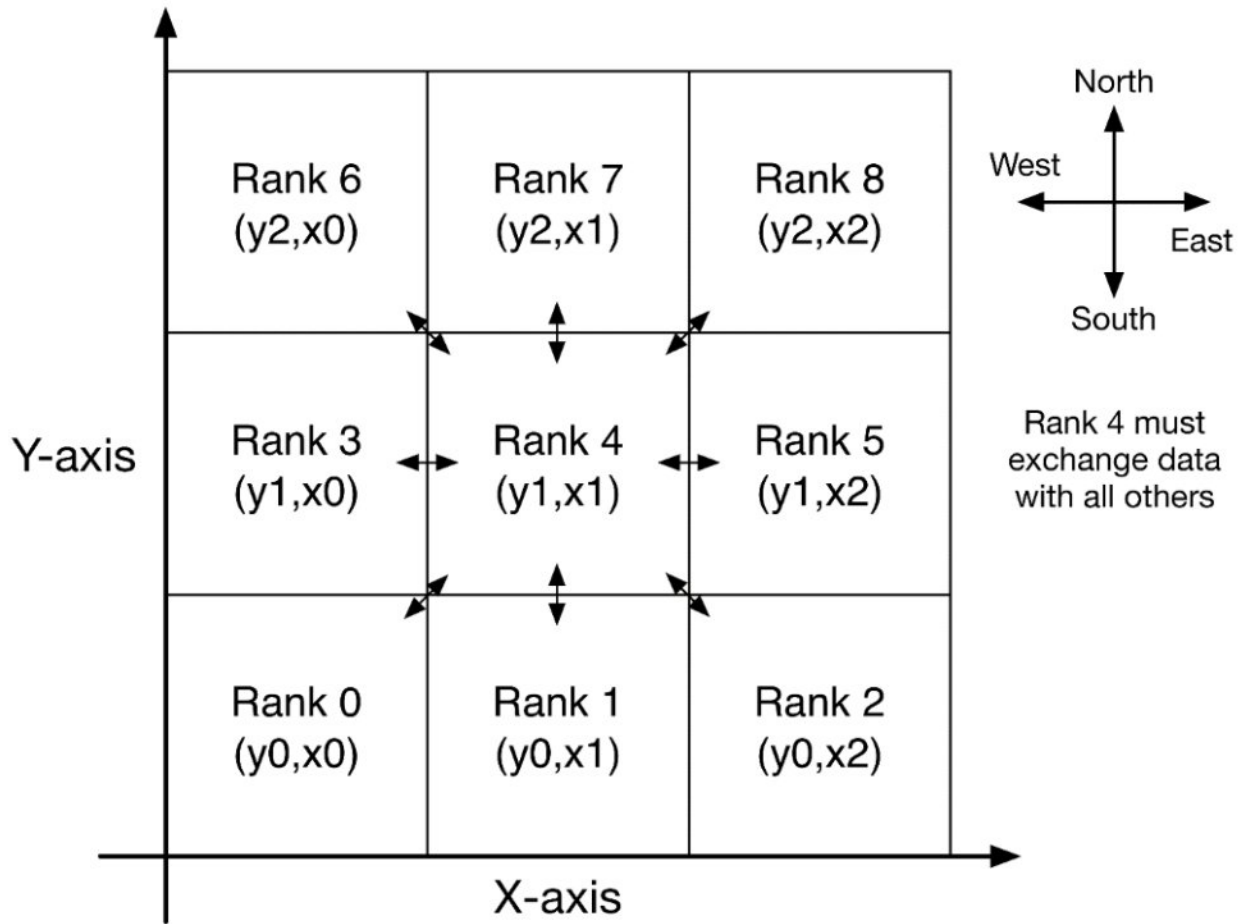
CPU Memory Energy Disk Network

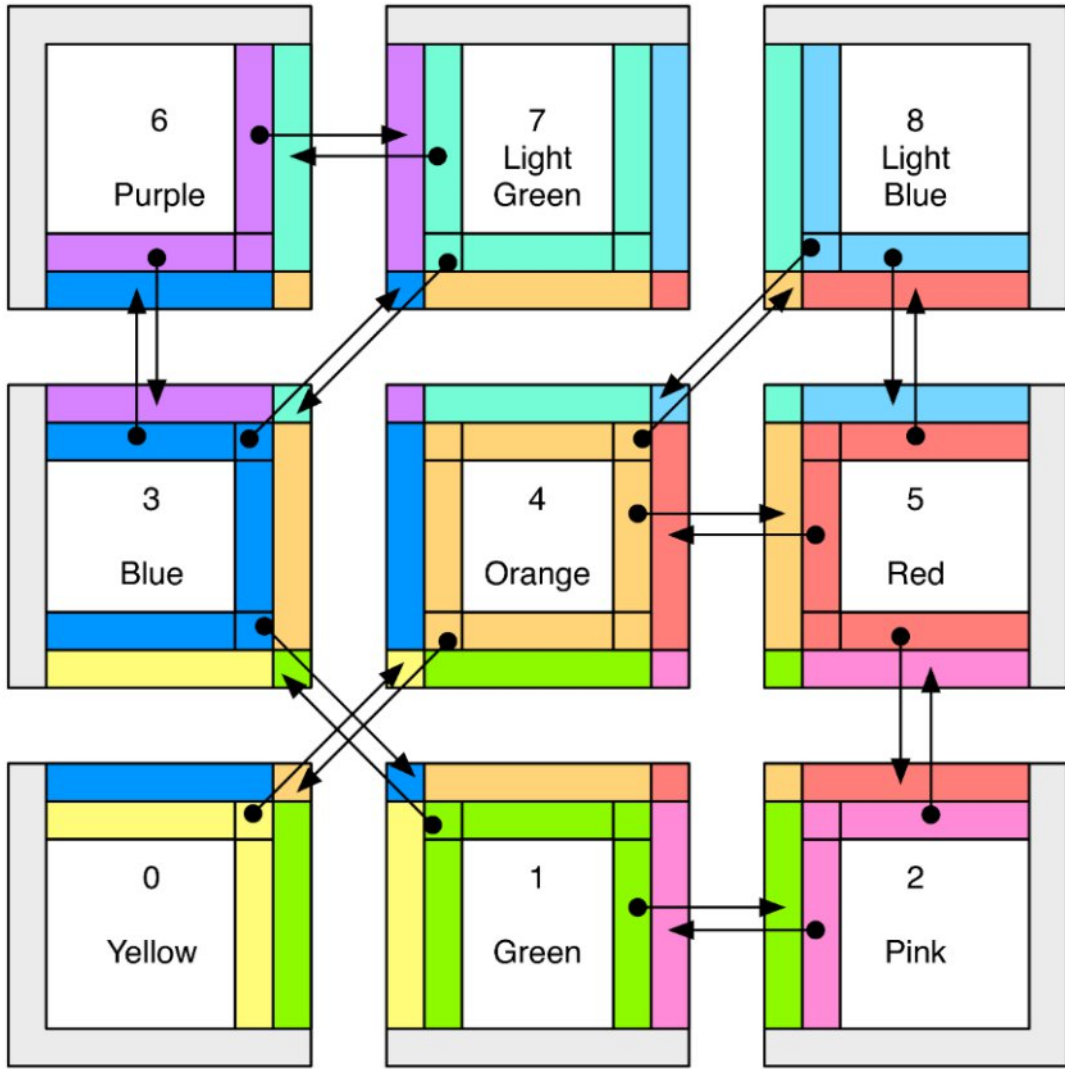
Process Name	Sent Bytes	Rcvd B...	Sent Pac...	Rcvd Packets	PID	User
R	3 KB	9 KB	44	54	12780	simon
R	3 KB	9 KB	23	32	12782	simon
R	3 KB	9 KB	23	34	12785	simon
R	3 KB	8 KB	21	28	12784	simon



Chapter 3: Advanced Message Passing







100	102	103	104	103
100	101	102	102	102
98	99	121	98	99
100	100	101	103	103
100	101	102	104	104

The median replacement value for pixel 121 is 101

98, 99, 100, 101,
101,
102, 102, 103, 121

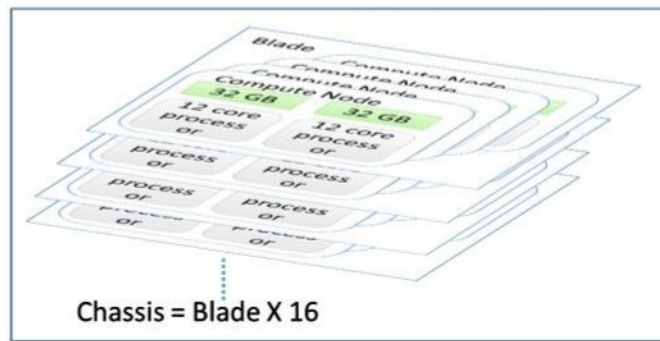
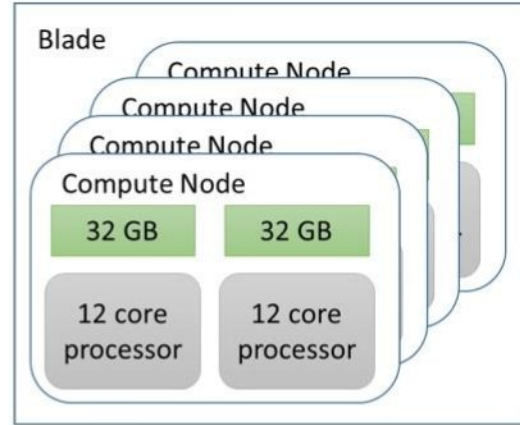
Chapter 4: Developing SPRINT, an MPI-Based R Package for Supercomputers

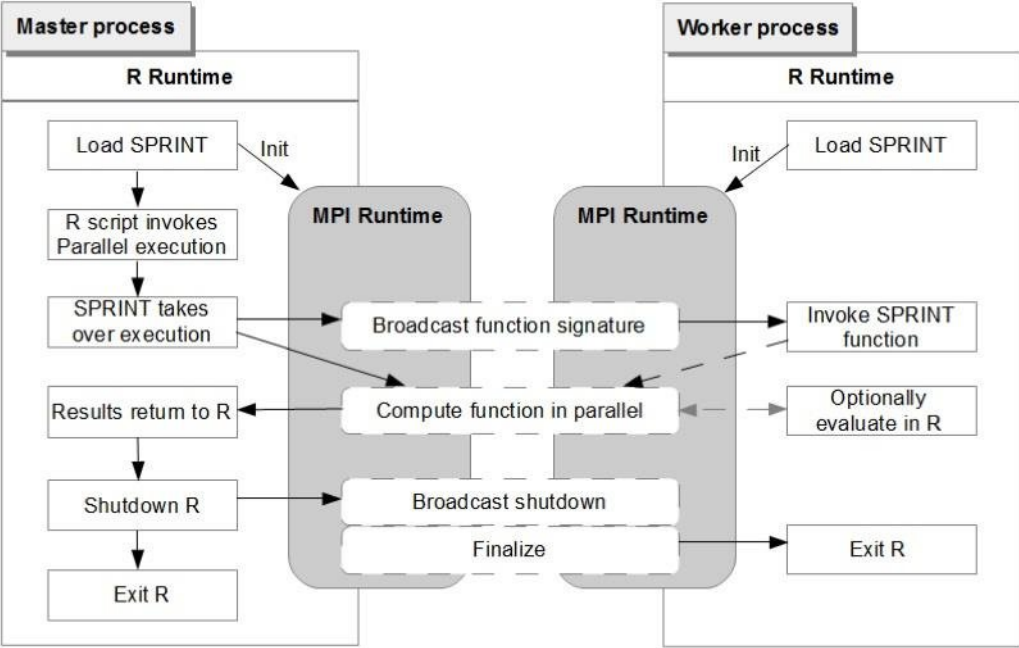




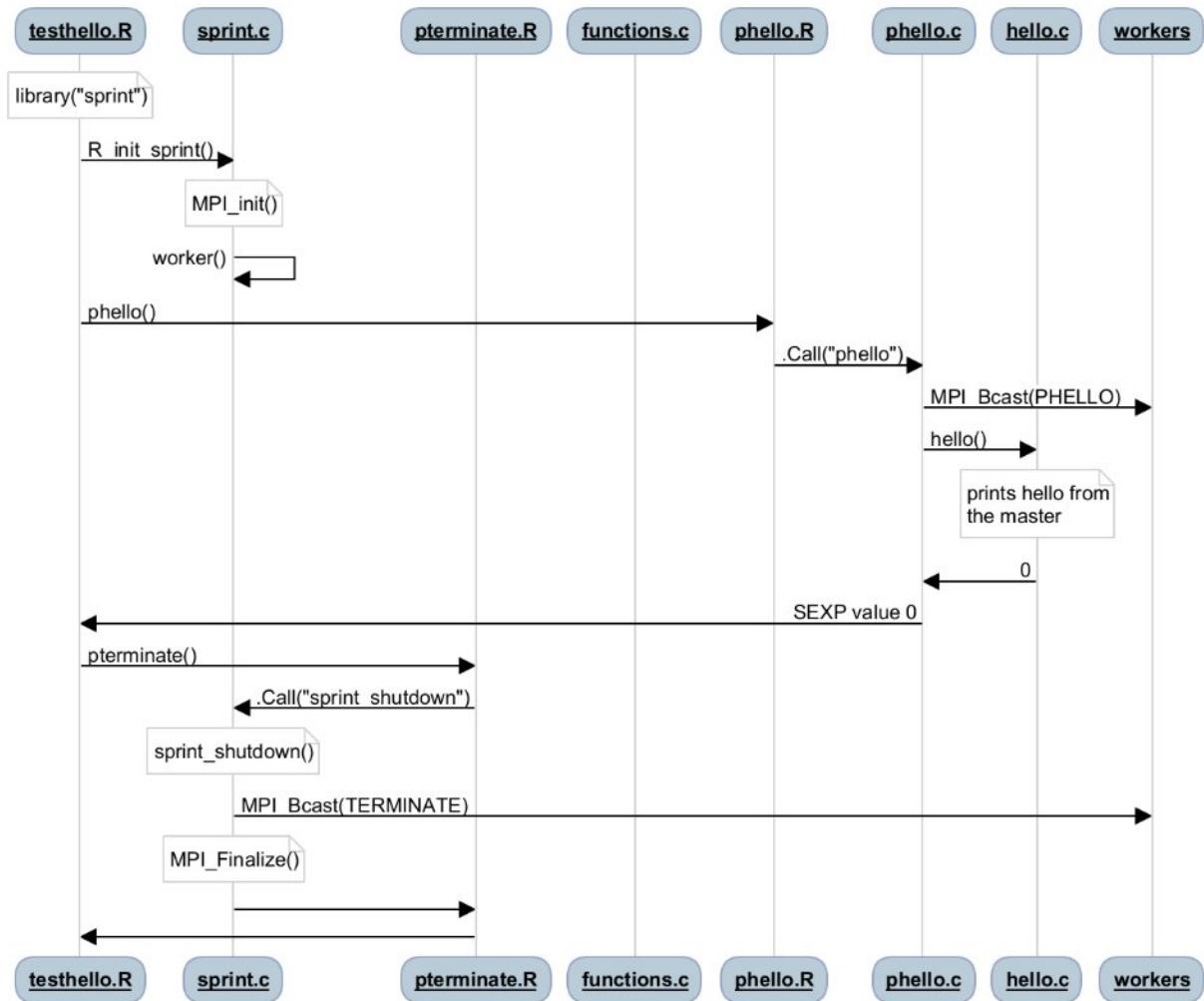
ARCHER has 4920 nodes in 26 cabinets giving a total of 118,080 cores

1 Cabinet
= 3 chassis

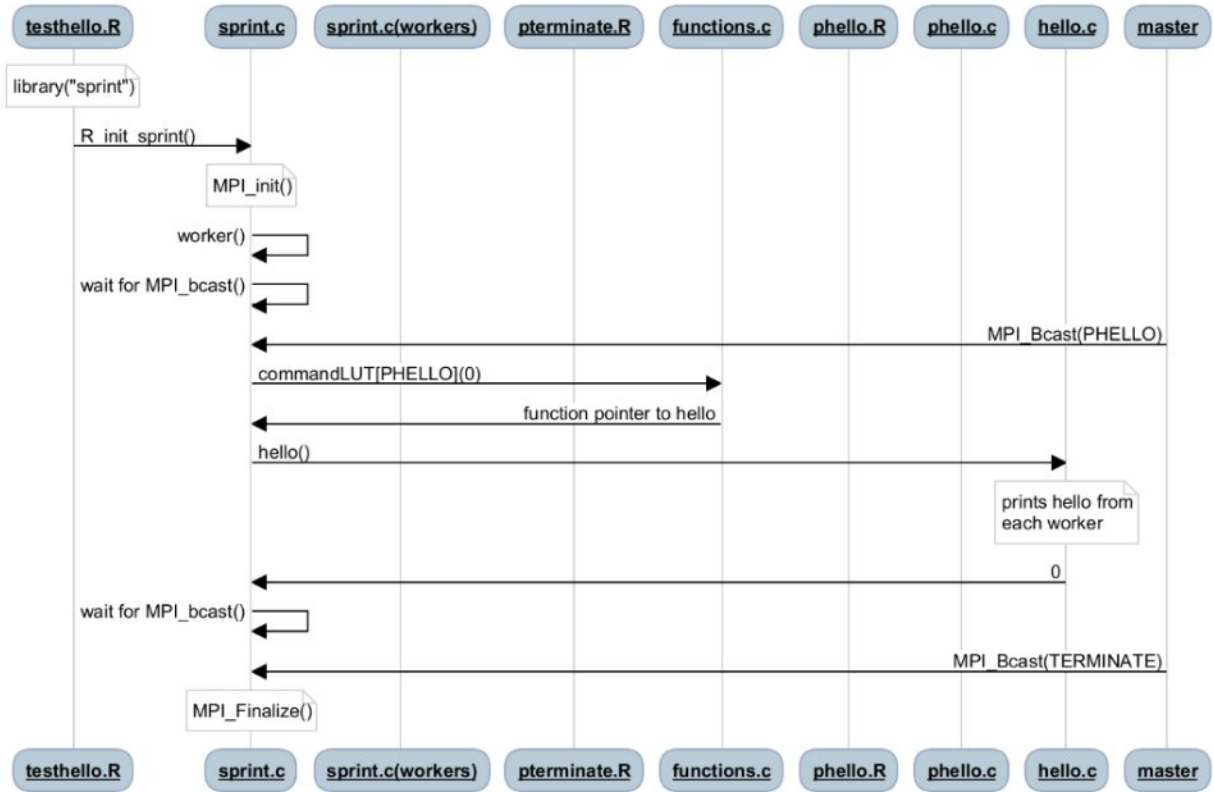


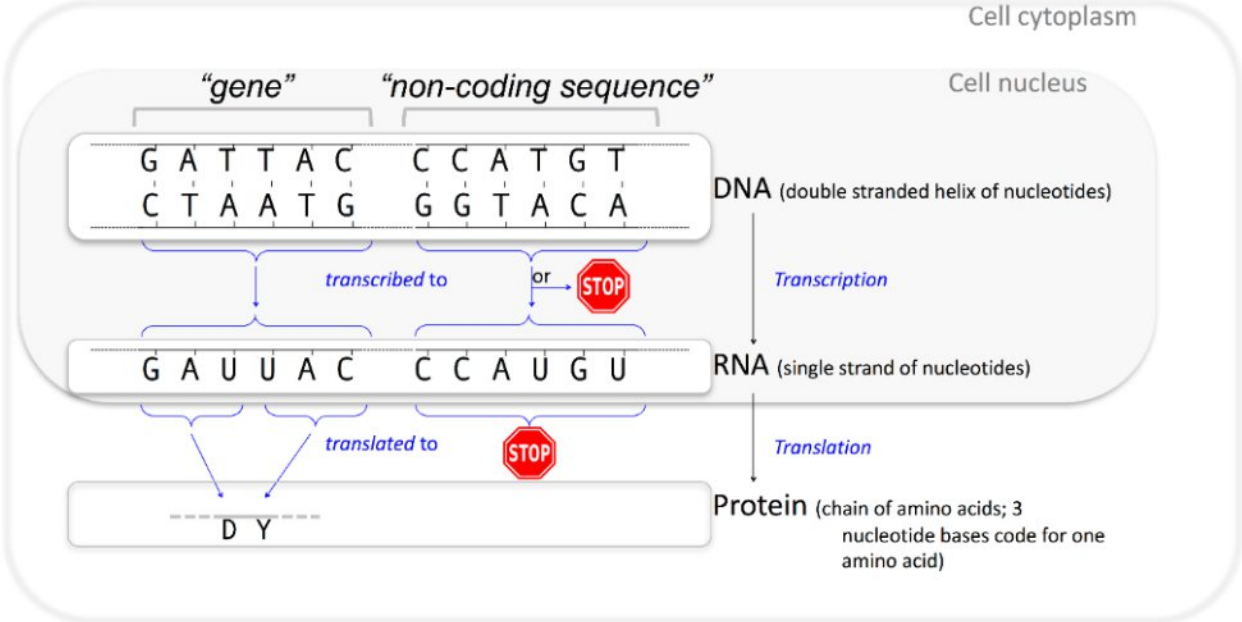


SPRINT hello world master

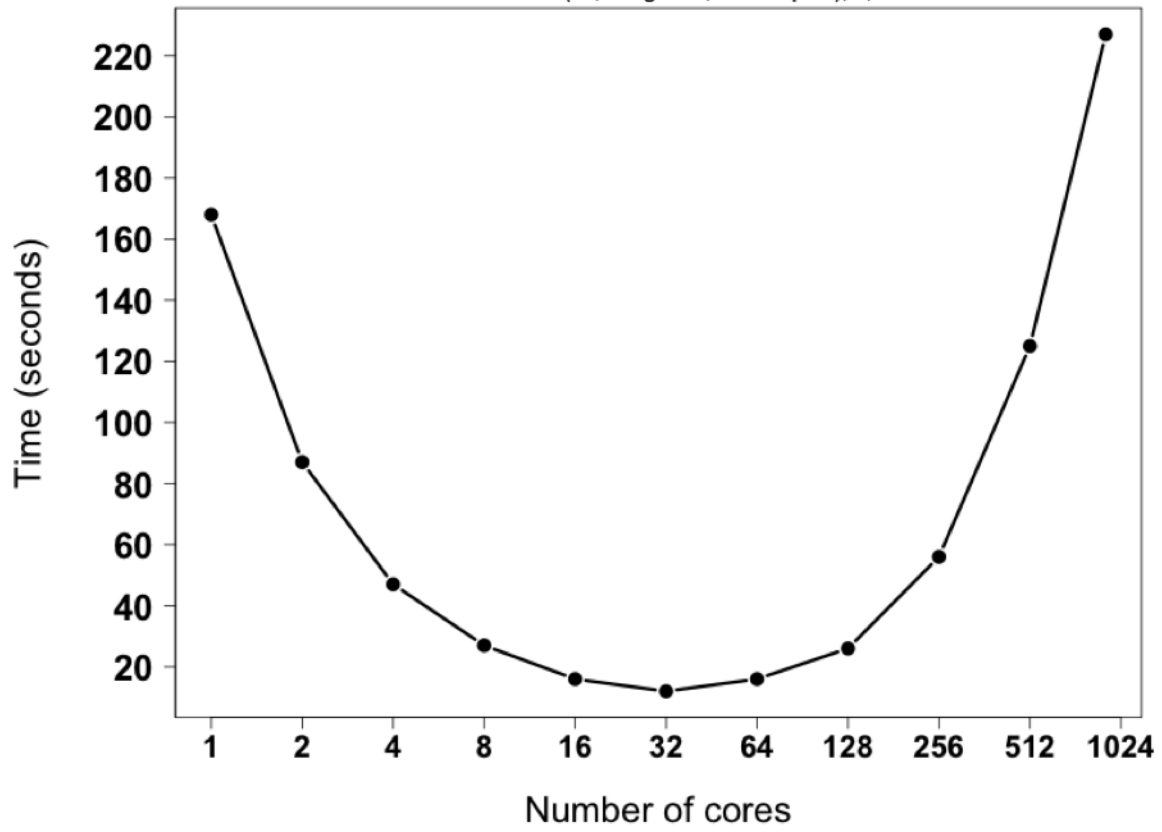


SPRINT hello world worker

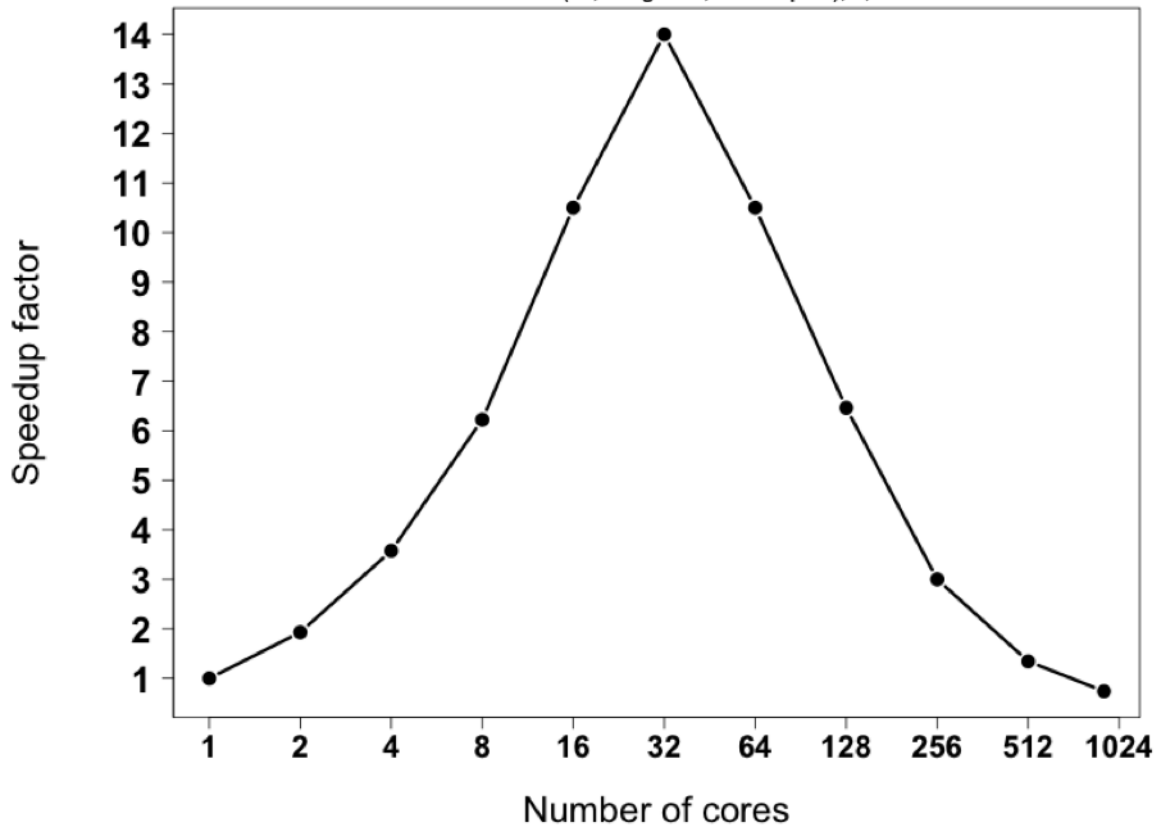




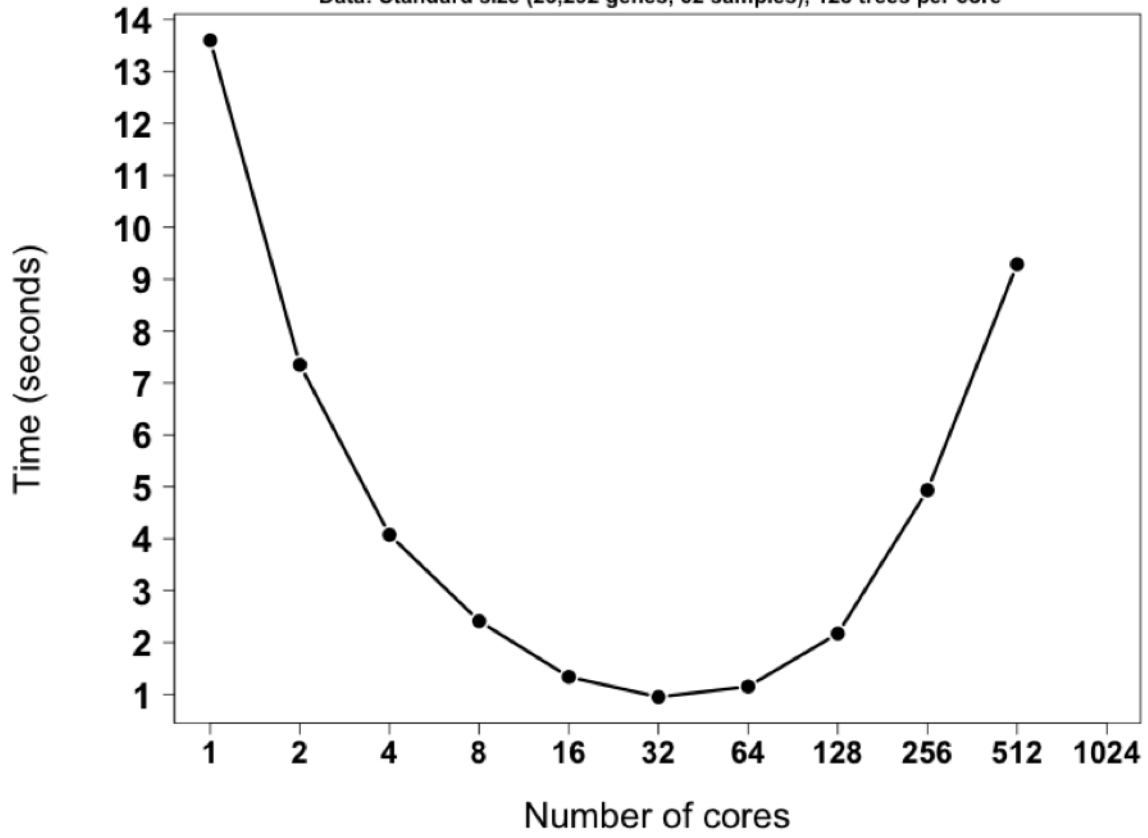
RandomForest run times
Data: Standard size (23,292 genes, 62 samples), 8,192 trees



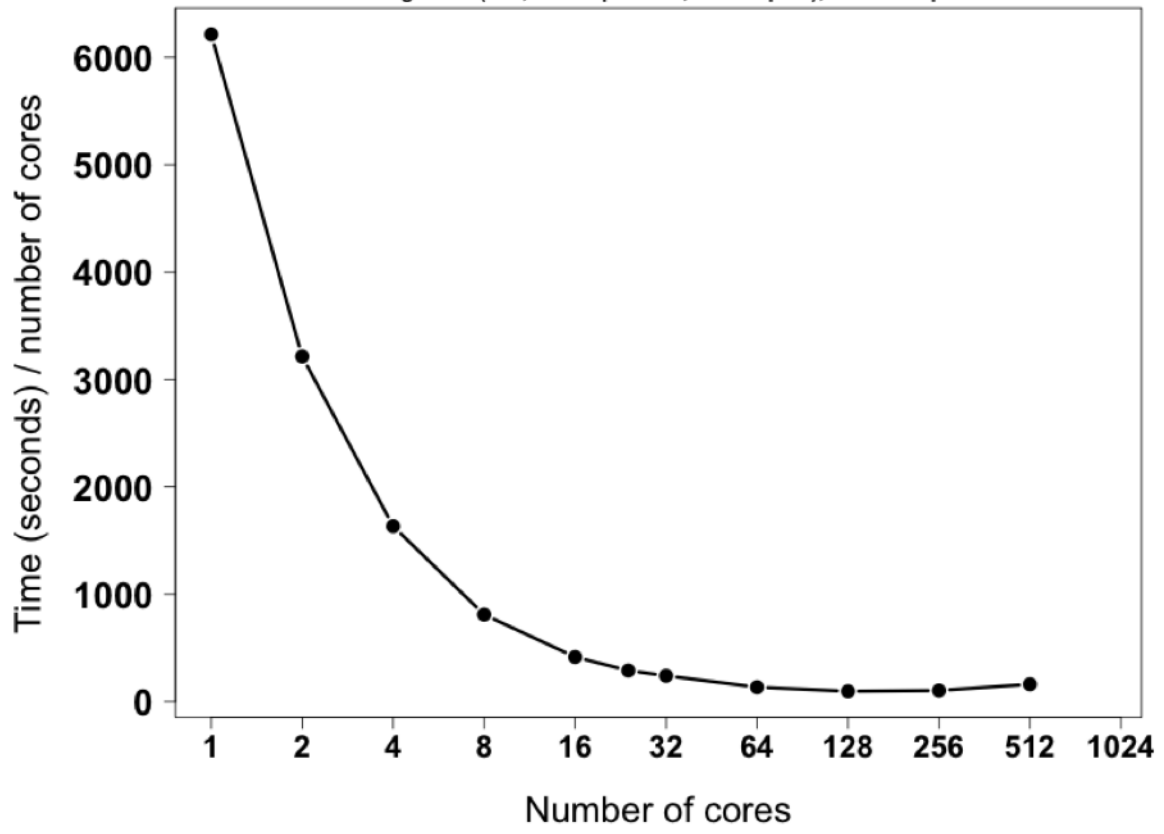
RandomForest speedup factors
Data: Standard size (23,292 genes, 62 samples), 8,192 trees



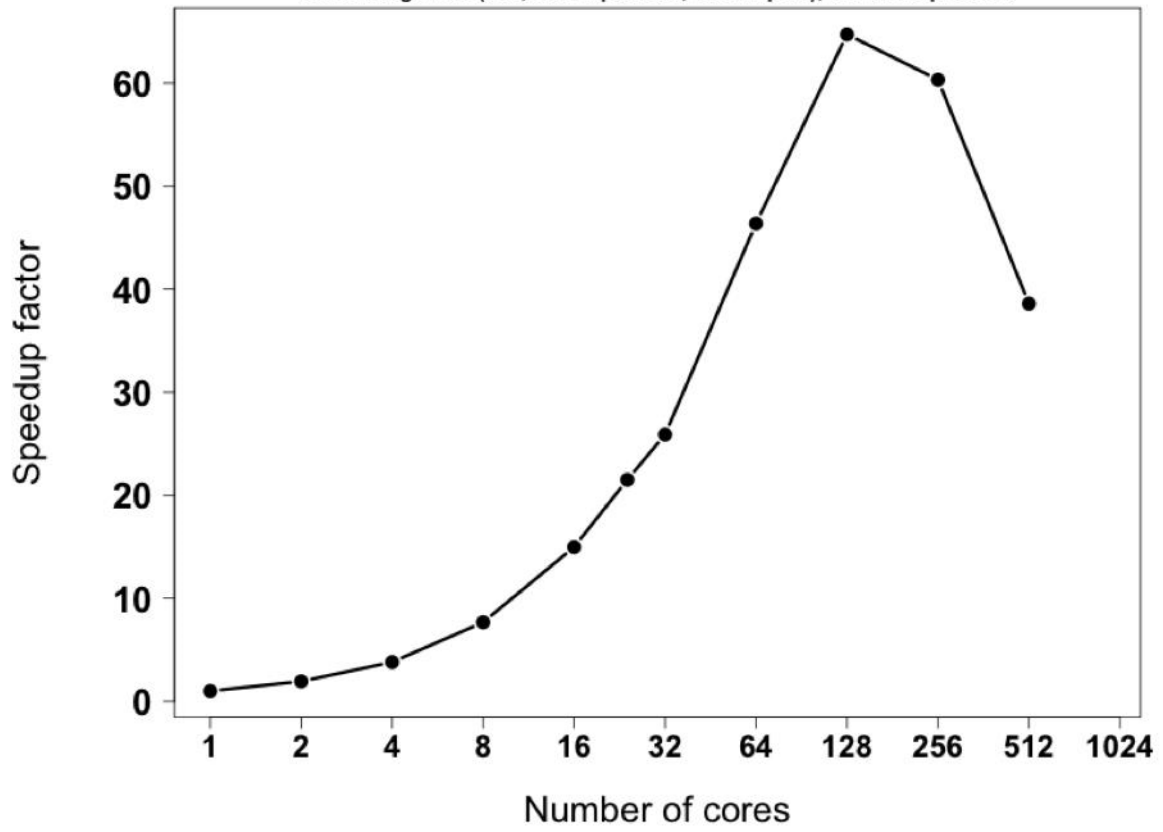
RandomForest speedup factors
Data: Standard size (23,292 genes, 62 samples), 128 trees per core



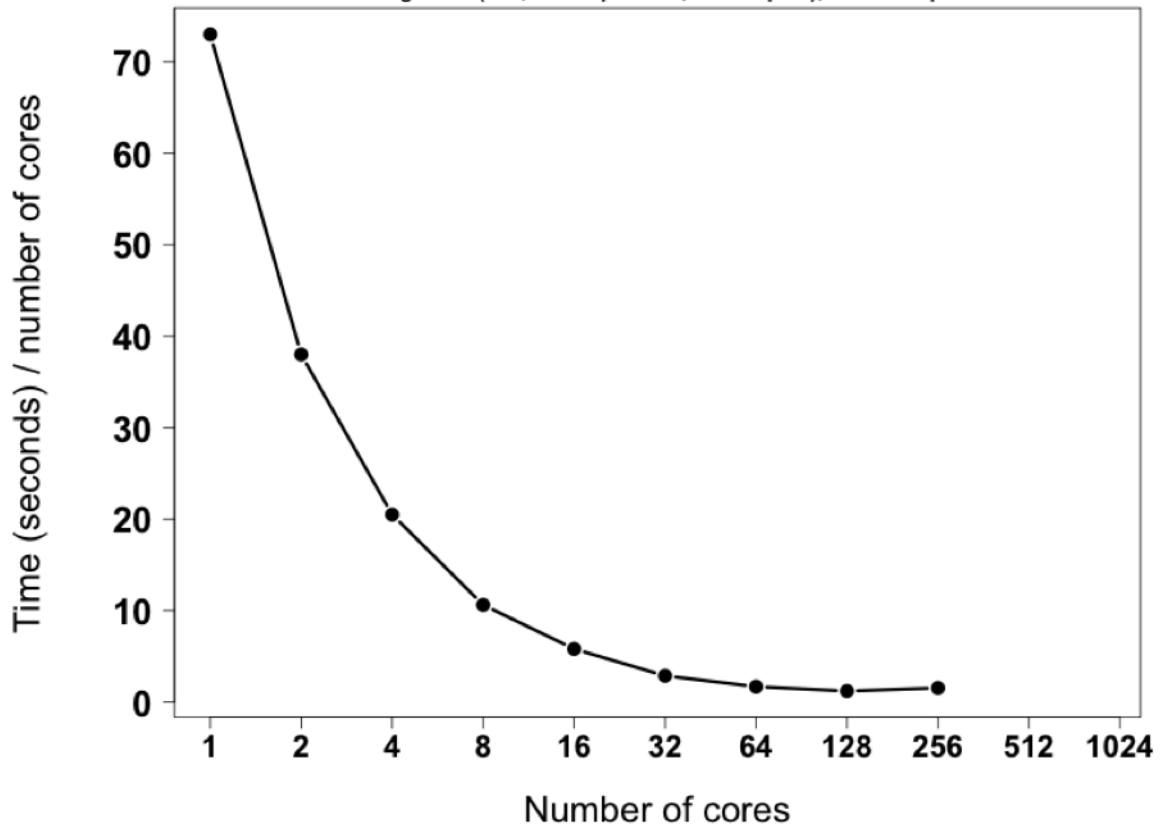
RandomForest run times
Data: Large size (500,000 sequences, 62 samples), 128 trees per core



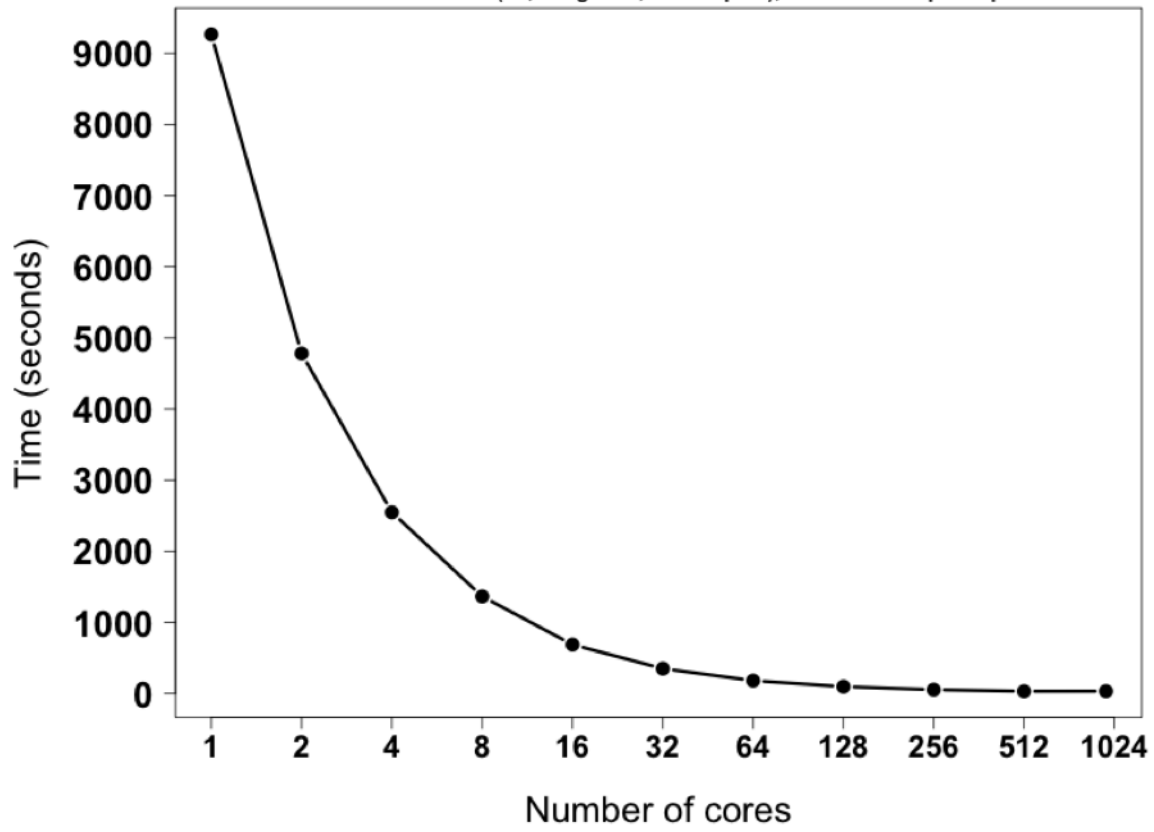
RandomForest speedup factors
Data: Large size (500,000 sequences, 62 samples), 128 trees per core

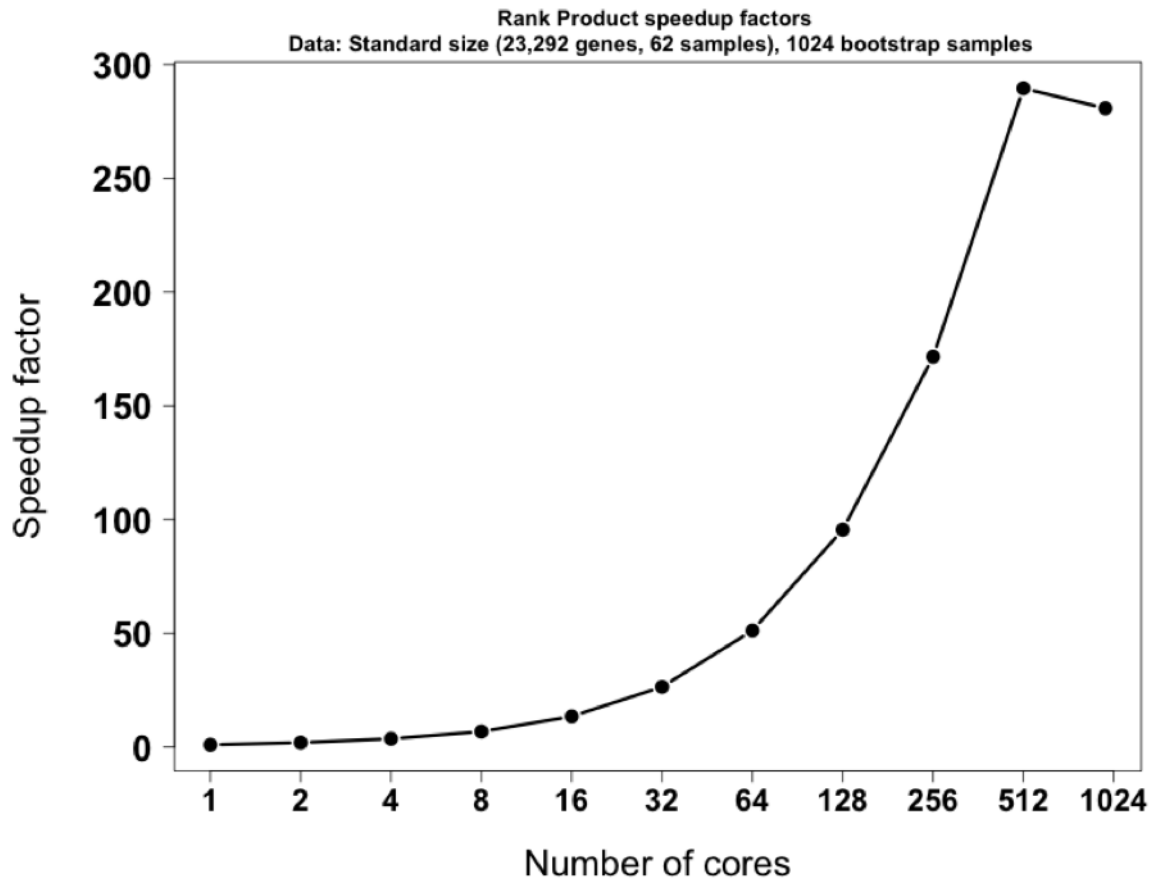


RandomForest run times
Data: Large size (500,000 sequences, 62 samples), 128 trees per core

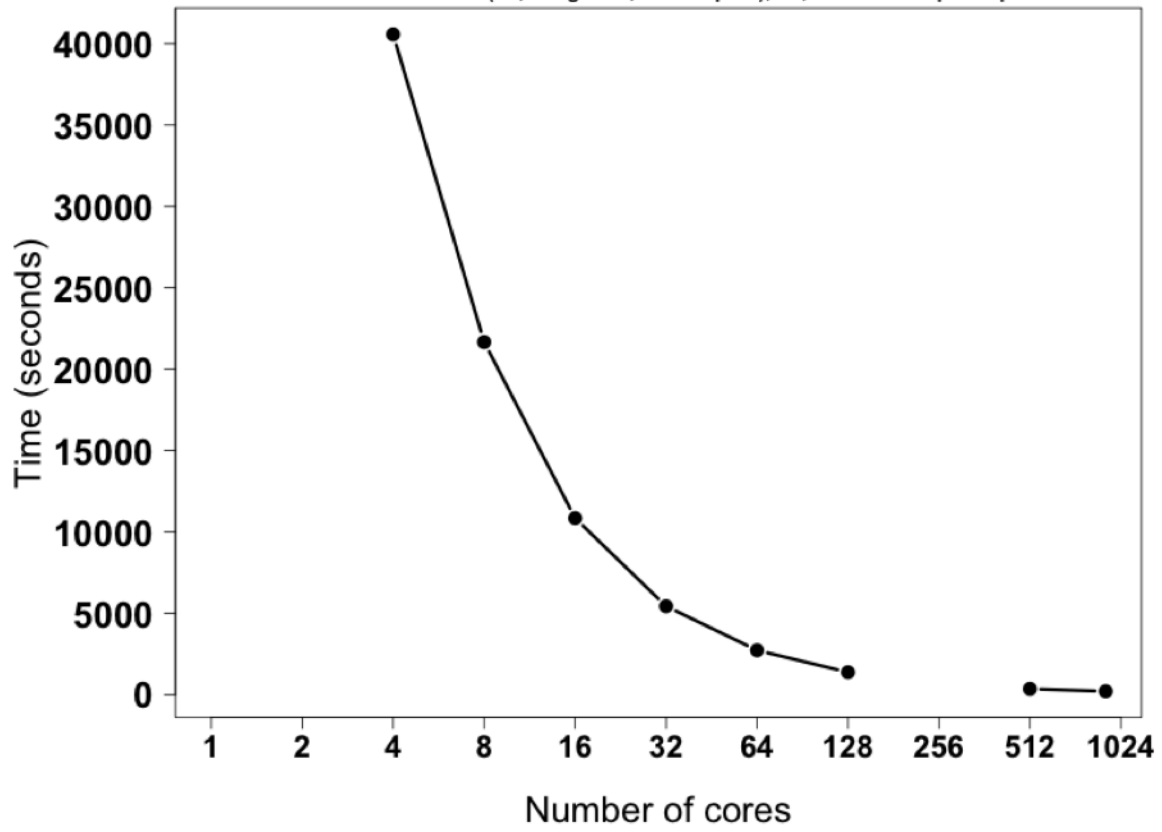


Rank Product run times
Data: Standard size (23,292 genes, 62 samples), 1024 bootstrap samples

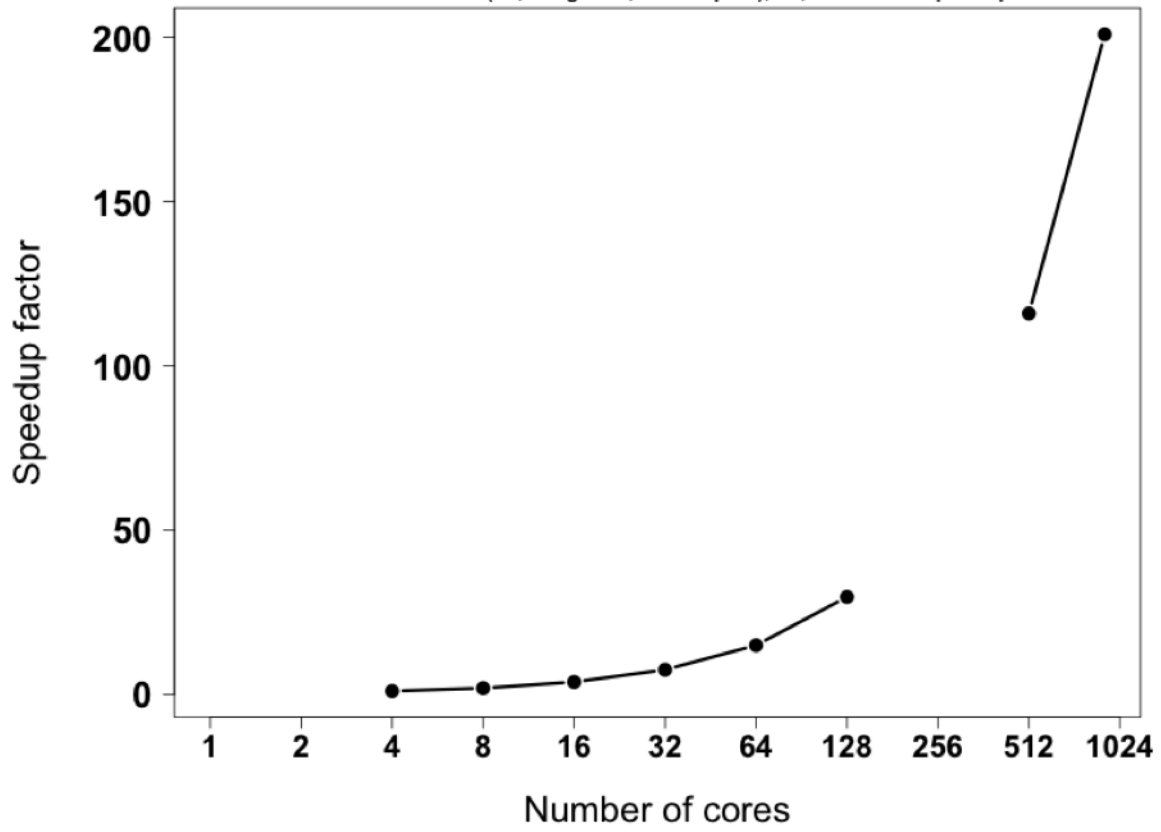




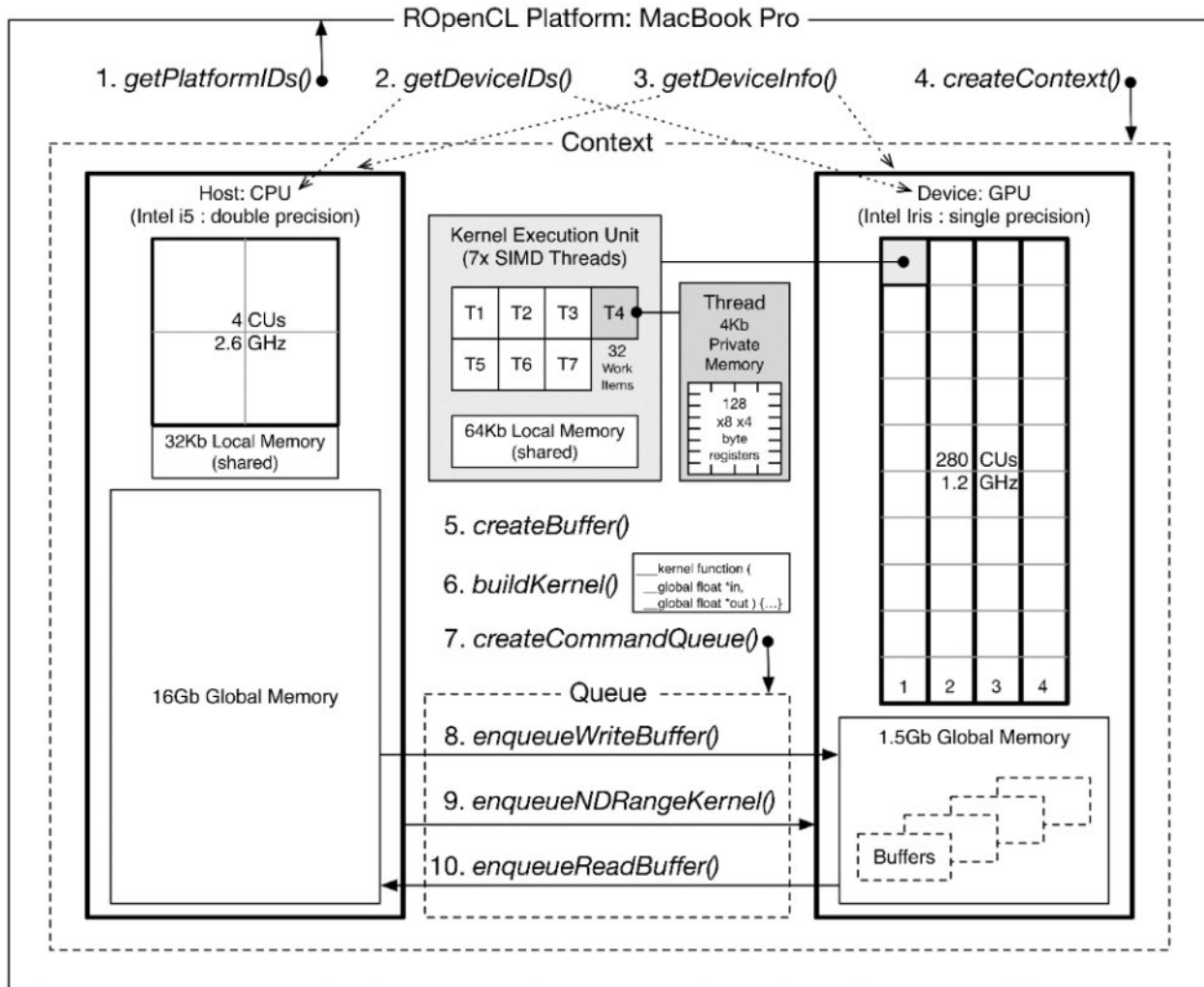
Rank Product run time
Data: Standard size (23,292 genes, 62 samples), 16,384 bootstrap samples



Rank Product speedup factors
Data: Standard size (23,292 genes, 62 samples), 16,384 bootstrap samples



Chapter 5: The Supercomputer in Your Laptop



$$\text{Euclidean distance} = \sqrt{\sum_{i=1}^N (A[i] - B[i])^2}$$

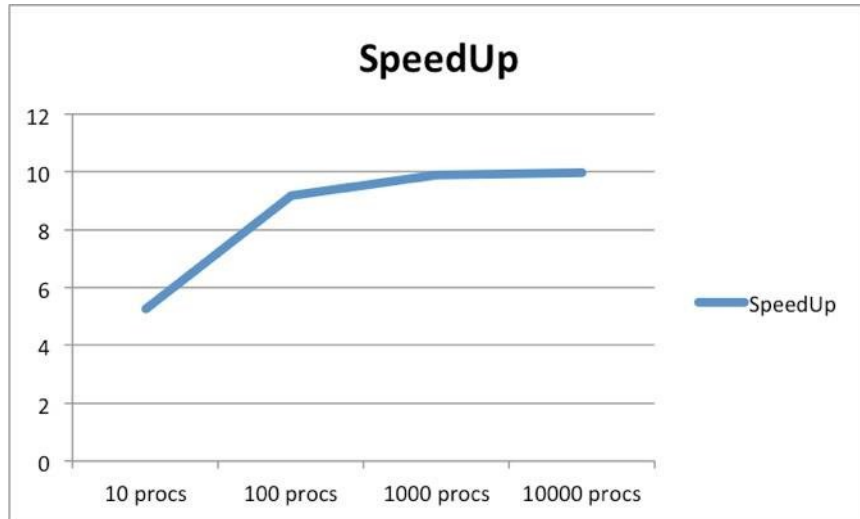
Chapter 6: The Art of Parallel Programming

$$\text{SpeedUp} = \frac{T_{\text{serial}}}{T_{\text{parallel}}}$$

$$\text{Perfect parallelism: } T_{\text{parallel}_N} = \frac{T_{\text{parallel}_1}}{N}$$

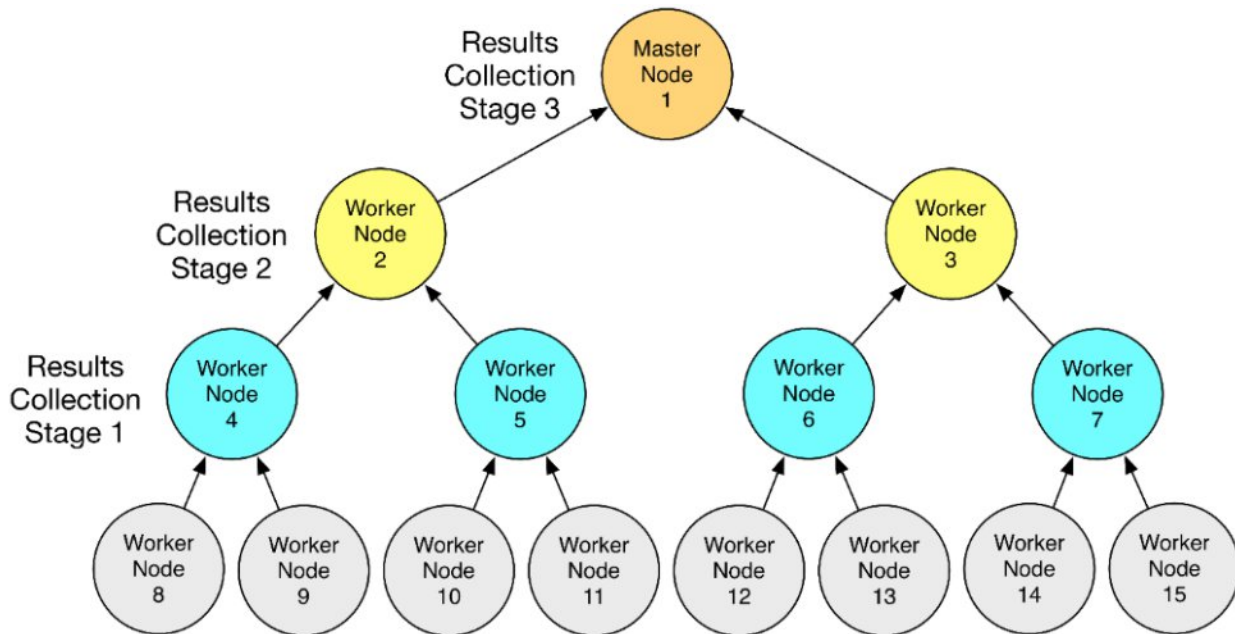
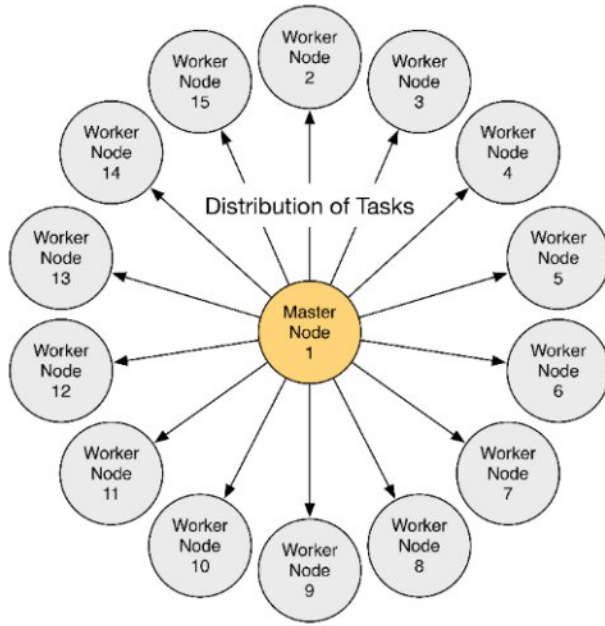
$$\text{Total time: } T_{\text{overall}_N} = T_{\text{non-parallel}} + T_{\text{parallel}_N}$$

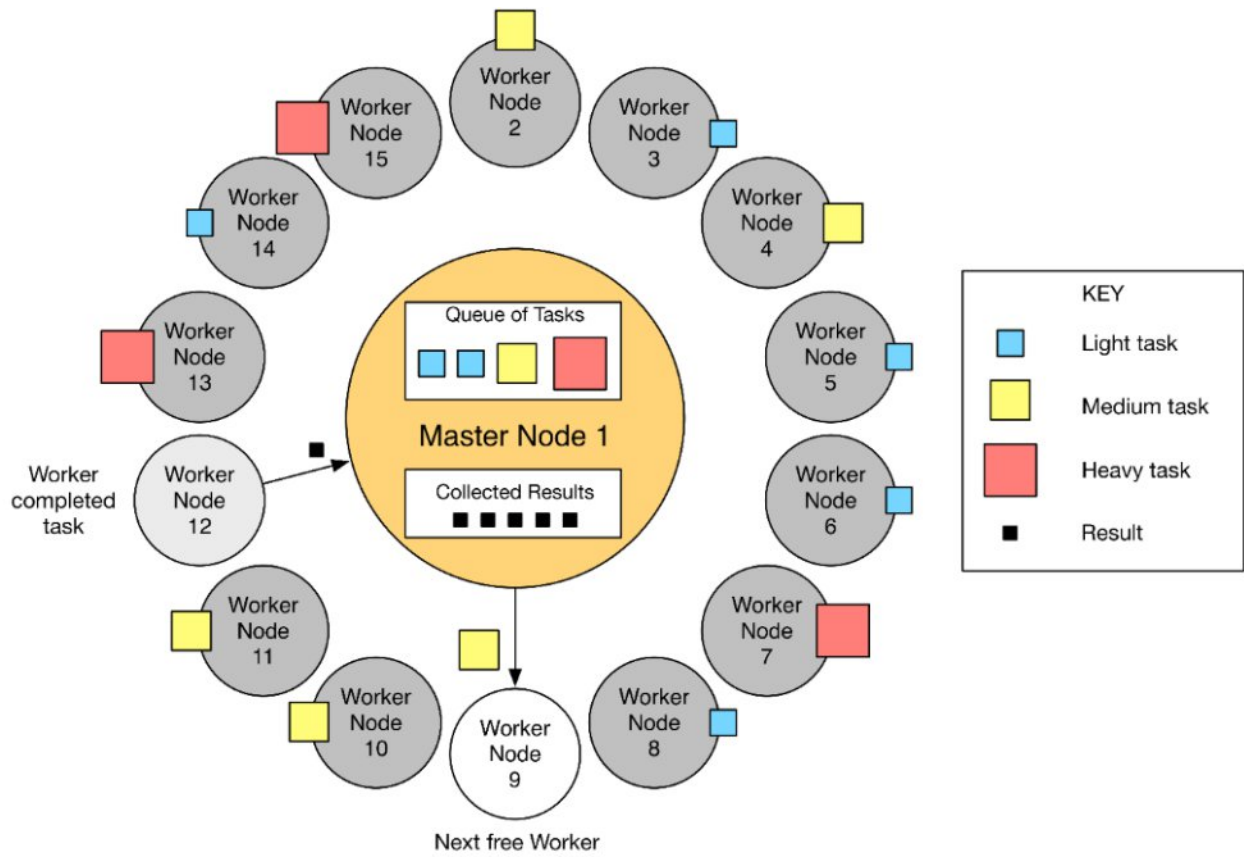
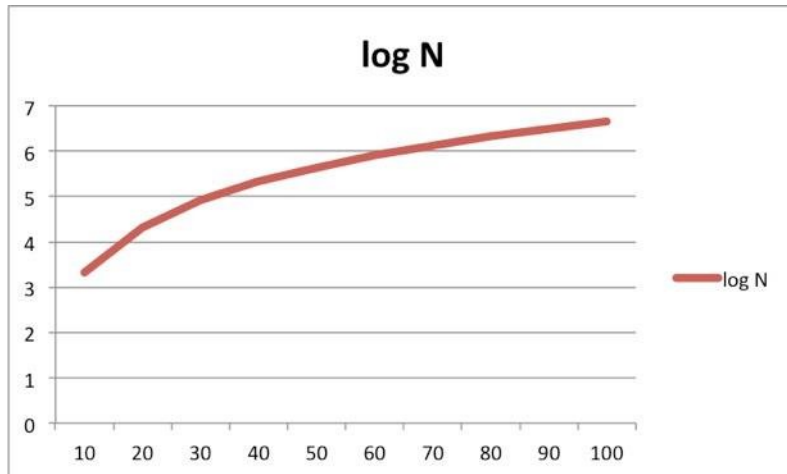
$$\text{SpeedUp}(N) = \frac{1}{(1-P) + \frac{P}{N}}$$

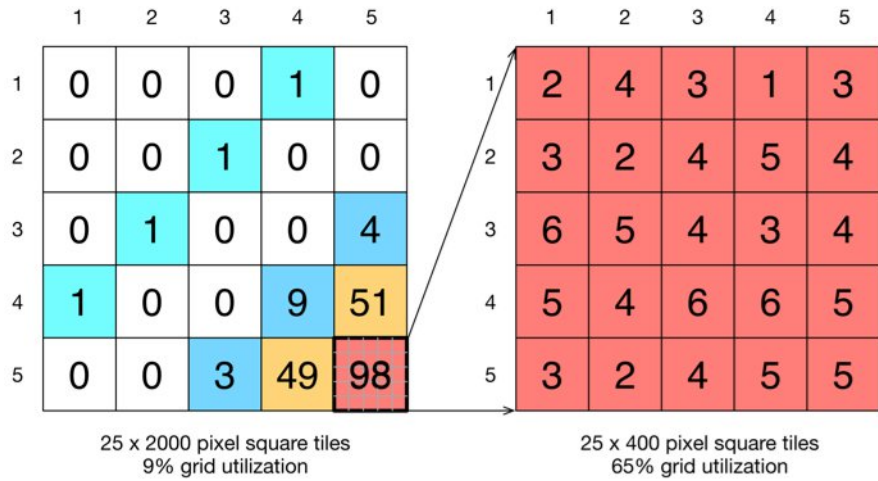
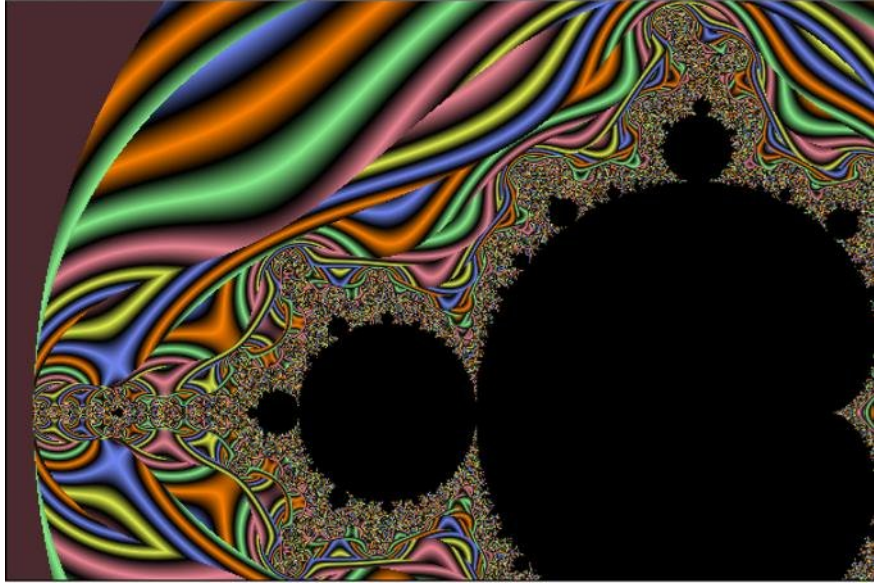


$$P_{estimated} = \frac{\frac{1}{SpeedUp} - 1}{\frac{1}{N} - 1}$$

Chapple's Law: $(T_{parallel} \times N) + T_{parallel_algorithm_development} \ll T_{serial} \times N$







Parallel computing: Applications

- The [caret](#) package by Kuhn can use various frameworks (MPI, NWS etc) to parallelized cross-validation and bootstrap characterizations of predictive models.
- The [maanova](#) package on Bioconductor by Wu can use [snow](#) and [Rmpi](#) for the analysis of micro-array experiments.
- The [pvclust](#) package by Suzuki and Shimodaira can use [snow](#) and [Rmpi](#) for hierarchical clustering via multiscale bootstraps.
- The [tm](#) package by Feinerer can use [snow](#) and [Rmpi](#) for parallelized text mining.
- The [varSelRF](#) package by Diaz-Uriarte can use [snow](#) and [Rmpi](#) for parallelized use of variable selection via random forests.
- The [bcp](#) package by Erdman and Emerson for the Bayesian analysis of change points can use [foreach](#) for parallelized operations.
- The [multtest](#) package by Pollard et al. on Bioconductor can use [snow](#), [Rmpi](#) or `rpvm` for resampling-based testing of multiple hypothesis.
- The [GAMBoost](#) package by Binder for `glm` and `gam` model fitting via boosting using b-splines, the [Geneland](#) package by Estoup, Guillot and Santos for structure detection from multilocus genetic data, the [Matching](#) package by Sekhon for multivariate and propensity score matching, the [STAR](#) package by Pouzat for spike train analysis, the [bnlearn](#) package by Scutari for bayesian network structure learning, the [latentnet](#) package by Krivitsky and Handcock for latent position and cluster models, the [lga](#) package by Harrington for linear grouping analysis, the [pepper](#) package by Porzelius and Binder for parallelised estimation of prediction error, the [orloca](#) package by Fernandez-Palacin and Munoz-Marquez for operations research locational analysis, the [rgenoud](#) package by Mebane and Sekhon for genetic optimization using derivatives the [afvyPara](#) package by Schmidberger, Vicedo and Mansmann for parallel normalization of Affymetrix microarrays, and the [puma](#) package by Pearson et al. which propagates uncertainty into standard microarray analyses such as differential expression all can use [snow](#) for parallelized operations using either one of the MPI, PVM, NWS or socket protocols supported by [snow](#).
- The [bugsparallel](#) package uses [Rmpi](#) for distributed computing of multiple MCMC chains using WinBUGS.
- The [partDSA](#) package uses [nws](#) for generating a piecewise constant estimation list of increasingly complex predictors based on an intensive and comprehensive search over the entire covariate space.
- The [dclone](#) package provides a global optimization approach and a variant of simulated annealing which exploits Bayesian MCMC tools to get MLE point estimates and standard errors using low level functions for implementing maximum likelihood estimating procedures for complex models using data cloning and Bayesian Markov chain Monte Carlo methods with support for JAGS, WinBUGS and OpenBUGS; parallel computing is supported via the [snow](#) package.
- The [pmclust](#) package utilizes unsupervised model-based clustering for high dimensional (ultra) large data. The package uses [pbdMPI](#) to perform a parallel version of the EM algorithm for finite mixture Gaussian models.
- The [harvestr](#) package provides helper functions for (reproducible) simulations.
- Nowadays, many packages can use the facilities offered by the **parallel** package. One example is [pls](#), another is [PGICA](#) which can run ICA analysis in parallel on SGE or multicore platforms.