

Analyzing Genetic Code Using Amazon Web Services

Catherine Vacher • David Wall

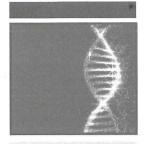
WILEY

สำนักหอสมุด มหาวิทยาลัยเชียงใหม่

6 16703959 0 1257676 i 22684396

Genomics in the AWS® Cloud





Genomics in the AWS® Cloud

Analyzing Genetic Code Using Amazon Web Services

> Catherine Vacher David Wall

> > WILEY

Contents at a Glance

Introduction		xix
Chapter 1	Why Do Genome Analysis Yourself When Commercial Offerings Exist?	•
Chapter 2	A Crash Course in Molecular Biology	9
Chapter 3	Obtaining Your Genome	25
Chapter 4	The Bioinformatics Workflow	39
Chapter 5	AWS Services for Genome Analysis	59
Chapter 6	Building Your Environment in the AWS Cloud	77
Chapter 7	Linux and AWS Command-Line Basics for Genomics	115
Chapter 8	Processing the Sequencing Data	143
Chapter 9	Visualizing the Genome	211
Chapter 10	Containerizing Your Workflow on the Desktop	235
Chapter 11	Variants and Applications	249
Chapter 12	Cancer Genomics	267
Index		291

Contents

Introduction)	xix
Chapter 1	Why Do Genome Analysis Yourself When Commercial Offerings Exist? Commercial Sequencing Services Typical Results Summary		1 2 3 8
Chapter 2	A Crash Course in Molecular Biology DNA DNA at Work: RNA and Proteins Inheritance Summary		9 13 20 23
Chapter 3	Obtaining Your Genome Preparing to Have Your Genome Sequenced Can It Affect My Insurance? Privacy Humility and Levelheadedness Validation with a Clinically Accredited Test Alternatives to Using Your Own Genome Specifying Lab Work		25 25 25 26 26 26 27 27
	Depth Sample Type Type of Output Files Sequencing Technology Genome vs. Exome vs. SNP Arrays Engaging a Laboratory Getting a Tissue Sample for DNA Extraction Rules and Regulations		27 28 28 28 30 30 31 32

	Do-It-Yourself Phlebotomy Legal Considerations Shipping the Sample Receiving the Results Sequences and Quality Control Information Alignment Information Variation Information Summary		33 34 35 36 36 37 38 38
Chapter 4	The Bioinformatics Workflow		39
_	Extraction of DNA		40
	Deriving Nucleated Cells from Whole Blood		40
	Processing Nucleated Cells		41
	FASTA Files	10	41
	FASTQ Files		42
	Phred Scores		44 44
	ASCII Encoding of Phred Scores		44
	Alignment to a Reference Genome Reference Genomes		48
	Quality Control		49
	Trimming		50
	The Alignment Process		51
	Marking Duplicates		53
	Recalibrating Base Quality Score		53
	Calling SNVs and Indel Variants		54
	Annotating SNVs and Indel Variants		55
	Prioritizing Variants		56
	Inheritance Analysis		56
	Identifying SVs and CNVs		57
	Bioinformatics Workflow		58
	Summary		58
Chapter 5	AWS Services for Genome Analysis		59
	General Concepts		61
	Networking		61
	AWS Functionalities		61
	AWS Accounts		61
	Virtual Private Cloud		62 63
	Subnets Elastic IP Addresses		65
	Custom Environments		65
	Storage		66
	S3		67
	Glacier		67
	Computing		68
	Elastic Compute Cloud		68
	Containers		70
	Lambda Functions		73

		Contents	XV .
	Workflow Management	74	
	AWS Batch	74	
	AWS Step Functions	74	
	Simple Workflow Service	75	
	Third-Party Solutions	75	
	Summary	75 75	
Chapter 6	Building Your Environment in the AWS Cloud	77	
	Setting Up a Virtual Private Cloud	77	
	Setting Up and Launching an EC2 Instance	82	
	Shutting Down an Instance to Save Money	91	
	Setting Up S3 Buckets	91	
	Configuring Your Account Securely	95	
	Turning On Multifactor Authentication	97	
	Establishing an AWS IAM Password Policy	101	
	Creating Groups	102	
	Creating Users	105	
	Setting Up Your Client Environment	106	
	Connecting to an EC2 Instance	106	
	Connecting from macOS or Unix/Linux	108	
	Connecting from Windows	109	
	Making S3 Buckets Available Locally	110	
	Mounting an S3 Bucket as a Windows Drive	111	
	T T	111	
	Mounting an S3 Bucket Under macOS and Linux Summary	113	
Chapter 7	Linux and AWS Command-Line Basics for Genomics	115	
	Selecting a Linux Distribution	115	
	Accessing Your AWS Linux Instance from Your		
	Local Computer	118	
	From Windows	118	
	From macOS	120	
	Options for Setting Up Linux on Your Personal Computer		
	Getting Familiar with the Command Line	123	
	Absolute and Relative References	124	
	Manipulating Files	126	
	Transferring Files to and from Your AWS Instance	127	
	Keyboard Shortcuts	128	
	Running Programs in the Background	128	
	Understanding File Permissions	129	
	Compressing and Archiving Files	130	
	Compression	131	
	Grep	132	
	Pipes and Redirection Operators	132	
	Text Processing Utilities: awk and sed	133	
	Managing Linux	135	
	Package Management Systems	135	
	i ackage ivialiagement bystems	133	

	The AWS Command-Line Interface Installing the AWS CLI Environment Windows macOS and Linux Configuring the AWS CLI Setting the Configuration at the Command Line Storing the Configuration in the Configuration File Testing Your Installation AWS CLI Essentials An Alternative Approach: AWS Systems Manager Summary	135 136 136 137 137 138 138 139 139 140 141
Chapter 8	Processing the Sequencing Data Getting from Data to Information Aligning to the Reference Genome Making Adjustments and Refinements to the Aligned Reads in the BAM File	143 143 145
	Identifying the Small Differences and Recording Them in the VCF File Making Adjustments and Refinements to the Variants in the VCF File Annotating the SNVs and Indels Prioritizing the Variants to Identify the Most Consequential Ones Trio Analysis and Inheritance Analysis Identifying and Annotating SVs and CNVs Setting Up AWS Services and Data Storage Copying the FASTQ Files Installing Docker and Containers Summary	155 157 160 162 164 167 172 196 197 210
Chapter 9	Visualizing the Genome Introducing Genome Visualizers Installing the IGV Desktop Visualizer Connecting the IGV Visualizer to Our AWS Data Loading Data into the IGV Visualizer Visualizing Aligned Sequencing Reads in IGV Have a CIGAR Analyzing Variants in IGV Summary	211 214 216 220 226 229 230 233
Chapter 10	Containerizing Your Workflow on the Desktop Introducing Containerization Understanding and Using Docker Installing Docker on Your Local Machine Downloading a Docker Image Viewing Available Docker Images Pupping a Docker Container Interactively	235 235 239 240 241 242

		Contents	xvii
	Removing a Docker Image More on Using the Docker Hub Containers for Genomics Work Summary	243 244 244 248	
Chapter 11	Variants and Applications Polygenic Risk Scores Genome-wide Association Studies Calculating a Polygenic Score Metagenomics AlphaFold Predicting Protein Structure from Protein Sequence—A 50-Year Puzzle Installing and Running AlphaFold Viewing and Comparing AlphaFold Results Summary	249 249 251 254 255 256 258 261 266	
Chapter 12	Cancer Genomics Somatic Genomes Cancer Oncogenes Tumor Suppressors The Promise and Reality of Cancer Precision Medicine Somatic or Germline? Cancer Predisposition Chromothripsis Epigenetics of Cancer Mechanisms of Cancer Samples Somatic Variant Analysis Copy Number Changes Measuring Tumor Genomic Instability Summary Notes	267 268 268 269 270 273 274 275 276 279 284 287 288 289	
Index		291	