

Décryptage de dynamiques épigénomiques au cours de la thymopoïèse et de la spermatogénèse en appliquant une méthodologie de recherche reproductible à des données de séquençage à haut débit

Guillaume Charbonnier

4 Octobre 2019

Jury

Dr Carl HERRMANN (GL)DKFZ (Heidelberg)Dr Marco Antonio MENDOZA (CR2)ISSB (Paris)Dr Sophie ROUSSEAUX (DR2)IAB (Grenoble)Dr Catherine NGUYEN (DR2)TAGC (Marseille)Dr Denis PUTHIER (MCU)TAGC (Marseille)Dr Salvatore SPICUGLIA (DR2)TAGC (Marseille)

Rapporteur Rapporteur Examinateur Examinateur Co-directeur de thèse Directeur de thèse







Epigenomic Dynamics

One genome, thousands phenotypes



Adapted from Ulirsch et al., 2019, Nature Genetics



Adapted from Dogan and Liu, 2018, Nature Plants



Adapted from Dogan and Liu, 2018, Nature Plants



Adapted from Dogan and Liu, 2018, Nature Plants



Adapted from Dogan and Liu, 2018, Nature Plants



Adapted from Dogan and Liu, 2018, Nature Plants

Genome Complexity: From structural organization to epigenetic regulation



Genome Complexity: From structural organization to epigenetic regulation



Genome Complexity: From structural organization to epigenetic regulation



Seeing the big epigenomic picture by mapping blood cells





Necker and TAGC: Blueprint collaboration and beyond







umr U 1090 TAGC theories and approaches of genomic complexity

Vahid Asnafi Agata Cieslak *et al.* Salvatore Spicuglia Guillaume Charbonnier Denis Puthier



Thymopoiesis

Thymopoiesis as a subset of hematopoiesis





Diagram adapted from Cancer UK Research / Wikimedia Commons



H&E stained section of human infant thymus taken from Darthmouth College



H&E stained section of human infant thymus taken from Darthmouth College



Diagram taken from Rothenberg et al., 2008, Nature Reviews Immunology



Available experimental approaches





Biological objectives



Thymopoiesis

A new reference epigenome of human early T cell differentiation

Epigenomic landscape summarized by chromatin segmentation



A consistent epigenomic landscape of hematopoiesis



MCA on collapsed chromatin states



Enrichment analysis of active enhancers highlights lineagespecific signatures



MSigDB Pathway top ontology terms

TCR signaling in CD8+ T cells -T cell receptor signaling pathway TCR signaling in CD4+ T cells Genes involved in Immune System Cytokine Signaling in Immune system Adaptive Immune System Platelet activation, signaling and aggregation Hemostasis -PDGFR-beta signaling pathway Fc gamma R-mediated phagocytosis Chemokine signaling pathway Innate Immune System BCR signaling pathway B cell receptor signaling pathway Signaling by the B Cell Receptor (BCR)

0 10 20 30 40 50

Chromatin segmentation meets expected epigenetics status for known cell-type specific genes



A consistent epigenomic landscape of thymopoiesis



A progressive loss of plasticity during differentiation



Chromatin segmentation meets expected epigenetics status for known cell-type specific genes



Thymopoiesis

Epigenomic dynamics of distal regulatory regions

DNA methylation dynamics of distal regions



DNA hypomethylation is a hallmark of distal regulatory regions...



mean DNA methylation ratio on region (0: unmethylated; 1: methylated)
DNA hypomethylation is a hallmark of distal regulatory regions...



mean DNA methylation ratio on region (0: unmethylated; 1: methylated)

DNA hypomethylation is a hallmark of distal regulatory regions mostly irrespective of their activation status in T cells



DNA hypomethylation is a hallmark of distal regulatory regions mostly irrespective of their activation status in T cells



DNA hypomethylation is a hallmark of distal regulatory regions mostly irrespective of their activation status in T cells



DNA hypomethylation is a hallmark of distal regulatory regions mostly irrespective of their activation status in T cells



Chromatin opening dynamics of distal regions in early T cell differentiation...



Dynamically open distal regions in early T cell differentiation are always hypomethylated



Constitutively open chromatin distal regions in early T cell differentiation exhibits H3K27ac dynamics...



Constitutively open chromatin distal regions in early T cell differentiation exhibits H3K27ac dynamics...



Constitutively open chromatin distal regions in early T cell differentiation exhibits H3K27ac dynamics associated with transcription of nearby genes



42

Differential expression analysis...



Differential expression analysis highlights genes associated with a dynamic putative enhancer



Thymopoiesis

Epigenetic regulation of TCRA locus

Epigenetic dynamics of TCRA locus



Activators binds TCRA enhancer in early T cell differentiation...





Activators binds TCRA enhancer in early T cell differentiation but the locus activates late



Looking for a potential candidate repressor

TLX homeodomain oncogenes mediate T cell maturation arrest in T-ALL via interaction with ETS1 and suppression of TCR α gene expression

Physiological αβ T cell development (DN4 to CD4+8+ stage)



Original discovery by Dadi, Spicuglia, Asnafi et al., 2012, Cancer Cell Figure taken from King, Ntziachristos & Aifantis, 2012, Cancer Cell

Looking for a potential candidate repressor in transcriptomic data...



Looking for a potential candidate repressor in transcriptomic data highlights HOXA family genes



C13 - TF Only [85]

Mature T-lymphoblastic leukemias provide another evidence for Homeobox family genes repressive action on $E\alpha$



HOXA5-9 genes repression validated by biological approaches from Necker collaborators

Mechanistic model for $E\alpha$ activation in T cell differentiation



Reproducibility

Articles containing "Reproducibility crisis" in Pubmed



Year

Defining reproduciblity



Definitions from Goodman et al., 2016, Science Translational Medicine

Main requirements for methods reproducibility

	Data	ode Results Software
Component	Requirements	Explanations
Data	Findable Accessible Interoperable Reusable	FAIR data principles
Software	Portable	Tools with permissive licence User-level OS-independent package management Automatic environment deployment
Analysis Code	Traceable	Tools and parameters for each step Documentation of order, input and output
	Automation	Minimal human action required to reproduce analysis

Additional requirements for massive genomic data handling

	Data	Code Results Software
Component	Requirements	Explanations
Data	Findable Accessible Interoperable Reusable	FAIR data principles
Software	Portable	Tools with permissive licence User-level OS-independent package management Automatic environment deployment
	Traceable	Tools and parameters for each step Documentation of order, input and output
Analysis Code	Automation	Minimal human action required to reproduce analysis
	Scalable	Straightforward and efficient parallelization of tasks Stoppable and resumable analyses Clusters and Clouds supports

Elements of solution for reproducibility in bioinformatics



git clone URL_for_workflow/repo.git
cd repo
snakemake --use-conda

Reproducibility

Beyond reproducibility





Integration of multiple workflows into one analysis



Integrating workflows using subworkflows...



Integrating workflows using subworkflows may leads to poor parallelization



Integrating workflows using merged workflow lead to optimal parallelization


Integrating workflows using merged workflow allows to reduce code redundancy



Integrating workflows using merged workflow allows to reduce code redundancy







Write a unique workflow to compute any possible analyses

Store a bioinformatician thesis project, a whole research career, or a project of any size with multiple collaborators inside a reproducible, scalable and lightweight workflow

Metaworkflow limits Snakemake freedom to bare minimum



All rules can and should be coerced to write their outputs in their own unique directory



snakemake C/B/A/sample1.ext C/B/A/sample2.ext

Easier development and benchmarking



snakemake C/B/A/example1.ext C/A/example1.ext

Easier development and benchmarking



snakemake C/B/A/smp1.ext C/A/smp1.ext Cx/B/A/smp1.ext
x can describe any combinations of parameters for each rule.
Traceable + Flexible + Small codebase

My thesis filetree, unique workflow DAG and codebase



pprox 60 to

 $\approx 1 \text{ month}/100 \text{ CPUs}$

150 generalized rules + 100 specific rules

Remaining \approx 3 mo codebase after rm -rf out and \approx 3 to of private input data.











Acknowledgments





THE BEST THESIS DEFENSE IS A GOOD THESIS OFFENSE.