# Design and Analysis of NGS Experiments

#### Nastasja Kreim & Anke Busch

#### Bioinformatics Core Facility Institute of Molecular Biology, Mainz





July 6th, 2020



# Experimental Design



B Mainz

# Why experimental design?

- to enable unbiased comparison between subjects, conditions, treatment groups
- to account for random variation
- to establish a relationship between cause and effect
- to disentangle biological variability from technical variability
- to enable the generalisation of findings



#### Hypothesis-driven research

- Is drug A better than drug B?
- Is there a genetic interaction between gene X and gene Y?
- Are transfected cells behaving differently than control cells?



#### Correlation does not mean Causation

Ice cream consumption

Coronary heart disease



#### Association does not mean Causation





#### Correlation does not mean Causation





# Variability

#### Modes of Variability

- biological variability between subjects /samples
- variability between conditions /groups
- technical variability (e.g. sample extraction, library preparation)
- $\Rightarrow$  we want to determine the variability between groups



Basic rules of experimental design

- Blocking for known confounding factors
- Randomisation for unknown confounding factors
- Replication to estimate the variability within a group/condition



B Mainz

# Blocking

Creation of homogeneous sample sets (for known confounding factors) with a varying factor of interest.



1B Mainz

#### Randomised Block Design





#### Randomised Block Design





Example: flowcell design for testing differential expression



Auer, Paul L., and R. W. Doerge. "Statistical design and analysis of RNA sequencing data." Genetics 185.2/acleular Biology (2010): 405-416.

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

#### Example: flowcell design for testing differential expression





Auer, Paul L., and R. W. Doerge. "Statistical design and analysis of RNA sequencing data." Genetics 185.2/olecular Biology (2010): 405-416.

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

# Example: Cages





Klaus, Bernd. "Statistical relevancerelevant statistics, part I." The EMBO journal 34.22 (2015): 2727-2730.

# Replication

• to estimate the effect size

• to estimate how precise the effect size estimates are (SD)



# Replication

- to estimate the effect size
- to estimate how precise the effect size estimates are (SD)
- $\Rightarrow$  to generalise findings



#### Replication



Nastasja Kreim & Anke Busch - IMB Mainz

**Design and Analysis of NGS Experiments** 

Experimental Design Raw Data Analysis & Mapping

Downstream Analysis & Mapping Downstream Analysis QC - Quality Control

#### Replication



Liu, Yuwen, Jie Zhou, and Kevin P. White. "RNA-seq differential expression studies: more sequence or more sequence of the sequence of the

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

Experimental Design Raw Data Analysis & Mapping

aw Data Analysis & Mapping Downstream Analysis QC - Quality Control

#### Replication





Liu, Yuwen, Jie Zhou, and Kevin P. White. "RNA-seq differential expression studies: more sequence or more se

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

#### Replicates

#### **Technical Replicates**

Technical replicates are replicates which have the same biological sample as origin and are processed and measured multiple times.

#### **Biological Replicates**

- in vivo: samples from different individuals
- in vitro: are there biological replicates?



#### In vitro: Cell culture





#### In vitro: Cell culture





#### In vitro: Cell culture





#### In vitro: Cell culture



 $\Rightarrow$  a little bit better. More variance than before through split up higher in the hierarchy.



#### In vitro: Cell culture





#### In vitro: Cell culture



Not ideal either maybe the best solution depending on the circumstances.



#### In vitro: Cell culture



Not ideal either maybe the best solution depending on the circumstances.  $\Rightarrow$  ideal would be to have cell cultures from different individuals of the same cell type.



# Control groups

- confirm validity of the experiment
- reference an experimental manipulation is compared to
- positive and negative controls should be included
- control groups should be as similar as possible to your experimental groups



# Summary I

- Define your hypothesis and formulate your expectations
- Use an appropriate control
- Replicate
- Block for known confounding factors
- Randomise for unknown confounding factors



Raw Data (FASTQ File Format) Mapping

#### After the experiment...

# Data analysis





Mainz

(c) Santiago Muñoz Prado

Raw Data (FASTQ File Format) Mapping

## Reminder: RNA-seq and ChIP-seq

#### RNA-seq

- sequence expressed (m)RNA
- Goal: find differences in expression / splicing

#### ChIP-seq

- sequence DNA bound to DNA-binding proteins
- Goal: find binding sites of DNA-binding protein



Raw Data (FASTQ File Format) Mapping



Raw Data (FASTQ File Format) Mapping



Raw Data (FASTQ File Format) Mapping

#### Short reads

@HWI-ST558:257:C4AJHACXX:1:1101:2005:2735 1:N:0:ATCACG TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG +

@CCFFFDFHHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI @HWI-ST558:257:C4AJHACXX:1:1101:2458:2678 1:N:0:ATCACG AGGTAAACAGACCATTGGATGGGAGATAGCAAGAACAATAGACTCCCTCAG +

:@?4ADDDFFDBDGFEBGF<<;A@F8<A4?FGEBBFFG<BB?@FG@D>?FB @HWI-ST558:257:C4AJHACXX:1:1101:3208:2718 1:N:0:ATCACG AGGAGGAGGAAGGTGATATCACTGCACAATTTTTCATCTGTTATGATCAAT +

@CCFFFDFDDHHDAACCBCHHIIGHIIIIIIIIGHHEGGFEFHGGEHHGH @HWI-ST558:257:C4AJHACXX:1:1101:3358:2699 1:N:0:ATCACG AGTGTGCCATAGAGCATGCTTGCTATTCCTACAACCCATCCTCTTCAAGCC +

===DBBDFDHBDFHIGIGIHIIEGEHGHGH@F@FHII;GGGHGGIGGIGII @HWI-ST558:257:C4AJHACXX:1:1101:3627:2685 1:N:0:ATCACG TGGACATATTTTGCATATGTTATCAACATTCATTCTCAGCCCCTTAATGCA +



Int

Raw Data (FASTQ File Format) Mapping

#### Short reads

@HWI-ST558:257:C4AJHACXX:1:1101:2005:2735 1:N:0:ATCACG TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG

 @CCFFFDFHHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI

 @HWI-ST558:257:C4AJHACXX:1:1101:2458:2678

 1:N:0:ATCACG

 AGGTAAACAGACCATTGGATGGGAGATAGCAAGAACAATAGACTCCCTCAG

 +

:@?4ADDDFFDBDGFEBGF<<;A@F8<A4?FGEBBFFG<BB?@FG@D>?FB @HWI-ST558:257:C4AJHACXX:1:1101:3208:2718 1:N:0:ATCACG AGGAGGAGGAAGGTGATATCACTGCACAATTTTTCATCTGTTATGATCAAT +

@CCFFFDFDDDHHDAACCBCHHIIGHIIIIIIIGHHEGGFEFHGGEHHGH @HWI-ST558:257:C4AJHACXX:1:1101:3358:2699 1:N:0:ATCACG AGTGTGCCATAGAGCATGCTTGCTATTCCTACAACCCATCCTCTTCAAGCC +

===DBBDFDHBDFHIGIGIHIIEGEHGHGH@F@FHII;GGGHGGIGGIGII @HWI-ST558:257:C4AJHACXX:1:1101:3627:2685 1:N:0:ATCACG TGGACATATTTTGCATATGTTATCAACATTCATTCTCAGCCCCTTAATGCA +



ZINE
Raw Data (FASTQ File Format) Mapping

### Short reads

@HWI-ST558:257:C4AJHACXX:1:1101:2005:2735 1:N:0:ATCACG TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG +

@CCFFFDFHHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI @HWI-ST558:257:C4AJHACXX:1:1101:2458:2678 1:N:0:ATCACG AGGTAAACAGACCATTGGATGGGAGATAGCAAGAACAATAGACTCCCTCAG +

:@?4ADDDFFDBDGFEBGF<<;A0F8<A4?FGEBBFFG<BB?@FG@D>?FB @HWI-ST558:257:C4AJHACXX:1:1101:3208:2718 1:N:0:ATCACG AGGAGGAGGAAGGTGATATCACTGCACAATTTTTCATCTGTTATGATCAAT +

@CCFFFDFDDDHHDAACCBCHHIIGHIIIIIIIGHHEGGFEFHGGEHHGH @HWI-ST558:257:C4AJHACXX:1:1101:3358:2699 1:N:0:ATCACG AGTGTGCCATAGAGCATGCTTGCTATTCCTACAACCCATCCTCTTCAAGCC +



INT

Raw Data (FASTQ File Format) Mapping

### Short reads

@HWI-ST558:257:C4AJHACXX:1:1101:2005:2735 1:N:0:ATCACG TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG +

@CCFFFDFHHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI @HWI-ST558:257:C4AJHACXX:1:1101:2458:2678 1:N:0:ATCACG AGGTAAACAGACCATTGGATGGGAGATAGCAAGAACAATAGACTCCCTCAG +

:@?4ADDDFFDBDGFEBGF<<;A@F8<A4?FGEBBFFG<BB?@FG@D>?FB @HWI-ST558:257:C4AJHACXX:1:1101:3208:2718 1:N:0:ATCACG AGGAGGAAGGTGATATCACTGCACAATTTTTCATCTGTTATGATCAAT

@CCFFFDFDDHHDAACCBCHHIIGHIIIIIIIIGHHEGGFEFHGGEHHGH @HWI-ST558:257:C4AJHACXX:1:1101:3358:2699 1:N:0:ATCACG AGTGTGCCATAGAGCATGCTTGCTATTCCTACAACCCATCCTCTTCAAGCC +

===DBBDFDHBDFHIGIGIHIIEGEHGHGH@F@FHII;GGGHGGIGGIGII @HWI-ST558:257:C4AJHACXX:1:1101:3627:2685 1:N:0:ATCACG TGGACATATTTTGCATATGTTATCAACATTCATTCTCAGCCCCTTAATGCA +

Molecular Biology

Int

Raw Data (FASTQ File Format) Mapping

# FASTQ format / base qualities

#### Single read

+

@HWI-ST558:257:C4AJHACXX:1:1101:2005:2735 1:N:0:ATCACG TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG

**@CCFFFDFHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI** 

header sequence (header2) qualities



Raw Data (FASTQ File Format) Mapping

# FASTQ format / base qualities

#### Single read

+

@HWI-ST558:257:C4AJHACXX:1:1101:2005:2735 1:N:0:ATCACG TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG header sequence (header2) qualities

@CCFFFDFHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI

#### Quality translation

```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJ
|||||| | | | |
012345......20...26...31.....40.
```



Experimental Design Raw Data Analysis & Mapping QC - Quality Control

!"#\$%&'()\*+,-./0123456789:;<=>?@ABCDEFGHIJ

Raw Data (FASTQ File Format)

# FASTQ format / base qualities

### Single read

Quality translation

+

111111

012345.

@HWI-ST558:257:C4AJHACXX:1:1101:2005:2735 1:N:0:ATCACG header TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG sequence (header2)

**@CCFFFDFHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI** 

# Phred qual. score Q $Q = -10 \log_{10} P$ or $P = 10^{\frac{-Q}{10}}$

qualities



Raw Data (FASTQ File Format) Mapping

# FASTQ format / base qualities

#### Single read

+

 @HWI-ST558:257:C4AJHACXX:1:1101:2005:2735
 1:N:0:ATCACG

 TAGCGGAACCAAGTGAGGAACTATGCCAGAGTCTATTACCATCTGTATCTG
 9

header sequence (header2) qualities

ar Biolom

#### @CCFFFDFHHHHHIIJJGIJJJHJJIGIIHGJIIGJJJJEGIJIGHIJI

Quality translation	Phred qual. score $Q$
!"#\$%&'()*+,/0123456789:;<=>?@ABCDEFGHIJ	$Q = -10 \log_{10} P$
01234520263140.	or $P=10^{rac{-Q}{10}}$

	Phred score $Q$	prob. of incorrect base call $P$	base call accuracy
U	10	0.1 = 1 in 10	90%
	20	0.01 = 1 in 100	99%
	30	0.001 = 1 in 1000	99.9%
	40	0.0001 = 1 in 10000	99.99%
	http://en.wikipedia.org/wiki/Phred_quality_scor		

Raw Data (FASTQ File Format) Mapping

# FASTQ format / base qualities

NextSeq	Q-score	binning
---------	---------	---------

Q-score bins	new Q-score
2 – 9	6
10 - 19	15
20 - 24	22
25 - 29	27
30 - 34	33
35 - 39	37
$\geq 40$	40



Mainz

Raw Data (FASTQ File Format) Mapping

# FASTQ format / base qualities

### NextSeq Q-score binning

Q-score bins	new Q-score	accuracy bins	assigned accuracy
2-9	6	36.90 - 87.41%	74.88%
10 - 19	15	90.00 - 98.74%	96.84%
20 - 24	22	99.00 - 99.60%	99.37%
25 - 29	27	99.68 - 99.87%	99.80%
30 - 34	33	99.90-99.96%	99.95%
35 — 39	37	99.97 - 99.99%	99.98%
$\ge 40$	40	$\geq$ 99.99%	99.99%



Mainz

Raw Data (FASTQ File Format) Mapping



Raw Data (FASTQ File Format) Mapping



Raw Data (FASTQ File Format) Mapping

# Read mapping

ACTGATCCCATTTCATTCAAAAAATCAA	AATAAACTCGC	СТАААТСАСАСААССАА	ACCTAAAACT
C	TTTTCGCCA	ATTAAACGAC	C AGGCA
CTTTCGACCA			ACACCUAUAC
AATTAA	TCAGACAAAC	CCT	ACAGGIAIAC
CAGACAAACCC	GAAACGI	TGAAGT TACCAG	AAGGCC
AAAATCAGA	S CS	ΔͲͲΔΔCΔΔCΔ	GCATTGAACAGAAAG
מאמעיי	TACACACA	ATTANCANCA	
ACCAGGTAAAGA	CATT		AGACAAA
АССААААААТТТА	ACAGTT	GAACACA	ACAGTTAGACACA
ACAGTTATAGCA		ACAGATAGACA	A
ACAG	GATTTGTGAGAG	2	λርͲλርλርλርλλር
ACATAGACGGCAACAA		ACACACAT	listing of
ACAGACGATTTAACACA			ACAG( Molecular Biology
Nastasja Kreim & Anke Busch –	IMB Mainz De	esign and Analysis of NGS Ex	periments

Raw Data (FASTQ File Format) Mapping

# Read mapping



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

Raw Data (FASTQ File Format) Mapping

sequencing: error-prone

CA

ar Biolom

SNPs / InDels

# Read mapping

### Alignment of reads to genomic positions

- millions of short reads
- human genome:  $\sim 3 * 10^9$  bp
- repetitive / low complexity regions

 $\Rightarrow$  complex task, computationally expensive

#### Tools

C7

- large number of specialized NGS read mappers
- non-splice-aware: Bowtie, BWA, ... (differences in accuracy, speed, memory)
- splice-aware: STAR, TopHat, ... (differences in accuracy, speed, memory)
- mapping parameters can have huge impact on mapping results

Raw Data (FASTQ File Format) Mapping





Raw Data (FASTQ File Format) Mapping





Raw Data (FASTQ File Format) Mapping





Raw Data (FASTQ File Format) Mapping





Raw Data (FASTQ File Format) Mapping





Raw Data (FASTQ File Format) Mapping

# Types of mapping



Nastasja Kreim & Anke Busch - IMB Mainz

**Design and Analysis of NGS Experiments** 

Raw Data (FASTQ File Format) Mapping





Raw Data (FASTQ File Format) Mapping





Raw Data (FASTQ File Format) Mapping



Raw Data (FASTQ File Format) Mapping

### Visualization: browser tracks





Raw Data (FASTQ File Format) Mapping



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses





RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses





RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses





RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses





RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Quantification of reads



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Quantification of reads



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

### RNAseq: differential expression analysis

#### Steps

- RNA quantification
- On normalization
- data modeling, distribution estimation
- visualization



Mainz

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

### **RNA** quantification




RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Normalization - why?

#### Why normalization? What to normalize for?

the number of counts is related to:

mRNA expression level (proportional)

but also:

- the sequencing depth
- the transcript length

Thus, we need within and between sample normalization.



NAINE

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses





RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses





RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Data modeling

#### Data modeling / distribution estimation

- read counts modeled as following a negative binomial distribution (=poisson-like distribution with variance not only depending on mean)
- for each gene: calculate p-value for differential expression:
  - model read counts within conditions
  - assumption: genes with similar expression have similar variance
  - compare variation within conditions to that between conditions
- correct p-value for multiple testing (FDR)



ZINE

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Tools and plots

#### Example results from DESeq2 package

(Love et al., Genome Biology, 2014)



- variance within/between groups especially high for low read counts
- DESeq2: moderation of fold changes



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Other downstream analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Other downstream analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Other downstream analyses



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Other downstream analyses

#### GO term analysis

- assign GO terms to each significantly differentially expressed gene
- look for enriched GO terms in sign. up- or down-regulated genes
- tools (e.g.):
  - clusterProfiler (Yu et al., OMICS, 2012.)
  - DAVID (Huang et al., Nature Protocols, 2009.)



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Quantification of reads



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Peak calling: general idea

#### Peak calling

- identification of regions (peaks) that are enriched in the ChIP sample relative to the control with statistical significance
- most common approach: sliding window, count reads per window
- **peak** = enriched relative to control:



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## ChIP vs. input control



#### ChIP

adapted from

http://saltmanquarterly.wordpress.com/2013/06/16/epigenetics-changing-how-we-interpret-genomes/

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## ChIP vs. input control



#### ChIP

#### input control



adapted from

http://saltmanquarterly.wordpress.com/2013/06/16/epigenetics-changing-how-we-interpret-genomes/

Nastasja Kreim & Anke Busch – IMB Mainz

Design and Analysis of NGS Experiments

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Recommended controls

#### Input (most popular)

#### cross-linked and sonicated (fragmented) DNA, but not IP'd



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Recommended controls

Input (most popular)

cross-linked and sonicated (fragmented) DNA, but not IP'd

#### IgG / mock IP

Immunoglobulin G (IgG) used as a control (unspecific) antibody, which does not recognize DNA or chromatin associated proteins



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Recommended controls

Input (most popular)

cross-linked and sonicated (fragmented) DNA, but not IP'd

#### IgG / mock IP

Immunoglobulin G (IgG) used as a control (unspecific) antibody, which does not recognize DNA or chromatin associated proteins

#### untagged epitope

in case of epitope tagged constructs, perform ChIP on cells lacking epitope tag



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Peak calling: general idea

#### Peak calling

- identification of regions (peaks) that are enriched in the ChIP sample relative to the control with statistical significance
- most common approach: sliding window, count reads per window
- **peak** = enriched relative to control:



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Peak calling: Strand specific profiles at enriched sites



- DNA sequences are sequenced from the 5' end
- alignment to genome results in two peaks (one on each strand)
- peaks are flanking the binding location of the protein of interest



RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Peak calling: construction of combined signal profiles



Nastasja Kreim & Anke Busch – IMB Mainz

**Design and Analysis of NGS Experiments** 

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

Peak calling: enrichment for TFs and histone modifications





Wilbanks & Facciotti, PLOS ONE, 2010

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Peak calling: enrichment for TFs and histone modifications





Wilbanks & Facciotti, PLOS ONE, 2010

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

## Peak calling: enrichment for TFs and histone modifications



Institute of

Wilbanks & Facciotti, PLOS ONE, 2010

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# Peak calling: tool comparison

#### Which peak finder should I use?

- dozens of different peak finders published
- some optimized for either TFs or histone marks
- sensitive to parameter settings
- e.g. MACS2 (https://github.com/macs3-project/MACS)

Reviews		
TF ChIP-seq:	histone ChIP-seq:	
• Laajala et al., BMC Genomics, 2009	<ul> <li>Micsinai et al., NAR, 2012</li> </ul>	
• Wilbanks & Facciotti, PLOS ONE, 2010		
General:		
• Pepke et al., Nat. Methods, 2009		nstitute of ar Biology

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# After peak calling: annotation / functional analysis

PeakAnalyzer (P	eakAnnotator)	(Salmon-Divon et al., BMC Bioinformatics, 2010)	
For each peak:	• downstream fo	downstream forward gene + distance	
	downstream re	o downstream reverse gene + distance	
	<ul> <li>overlapped ge overlap center</li> </ul>	• overlapped genes + overlap start (feat.) + overlap center (feat.) + overlap end (feat.)	
	Lico		
ChIPseeker	all	(https://github.com/GuangchuangYu/ChIPseeker)	
Peaks overlapping gene features			

- Other Intron (26.7%)
   Downstream (<=3kb) (1.5%)</li>
- Distal Intergenic (28.28%)

1st Intron (8,17%)

nstitute of ar Biology

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# After peak calling: differential binding sites



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

RNAseq: Differential Expression of Genes ChIPseq: Peak Calling ChIPseq: Annotation & Differential Binding Other subsequent analyses

# After peak calling: detection of (novel) binding motifs



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

Raw data quality control Mapping quality control Application specific quality control Batch effects



Raw data quality control Mapping quality control Application specific quality control Batch effects

# Quality control of next generation sequencing data



Raw data quality control Mapping quality control Application specific quality control Batch effects

## Library Preparation



Raw data quality control Mapping quality control Application specific quality control Batch effects





Raw data quality control Mapping quality control Application specific quality control Batch effects

# Raw data quality control

- quality score distribution
- base composition distribution
- read length distribution
- distribution of reads over samples
- overrepresented sequences



Raw data quality control Mapping quality control Application specific quality control Batch effects

# Distribution of quality values along reads



plot produced by fastqc: Andrews, S. "FASTQC. A quality control tool for high throughput sequence data. Molecular Biolog URL http://www. bioinformatics. babraham. ac. uk/projects/fastqc (2010).
Raw data quality control Mapping quality control Application specific quality contro Batch effects

### Quality score distribution



ata."Molecular Biology

plot produced by fastqc: Andrews, S. "FASTQC. A quality control tool for high throughput sequence data." Until the of URL http://www. bioinformatics. babraham. ac. uk/projects/fastqc (2010).

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Per base sequence content



plot produced by fastqc: Andrews, S. "FASTQC. A quality control tool for high throughput sequence data. Molecular Biology URL http://www. bioinformatics. babraham. ac. uk/projects/fastqc (2010).

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Per base sequence content





plot produced by fastqc: Andrews, S. "FASTQC. A quality control tool for high throughput sequence data." Molecular Biology URL http://www. bioinformatics. babraham. ac. uk/projects/fastqc (2010).

Raw data quality control Mapping quality control Application specific quality contro Batch effects

NB Mainz

#### Overrepresented sequences

#### Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GCCTAATTTAGGCAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTCGTA	834146	4.346867007222822	Illumina Paired End PCR Primer 2 (100% over 45bp)
GNCTAATTTAGGCAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTCGTA	515410	2.685883195738773	Illumina Paired End PCR Primer 2 (100% over 45bp)
GCCTAATTTAGGAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTCGTAT	27500	0.14330685839005114	Illumina Paired End PCR Primer 2 (100% over 46bp)
Bioinform			



plot produced by fastqc: Andrews, S. "FASTQC. A quality control tool for high throughput sequence data. Molecular Biolog URL http://www. bioinformatics. babraham. ac. uk/projects/fastqc (2010).









Raw data quality control Mapping quality control Application specific quality control Batch effects

# Alignment / Mapping quality control

- # reads mapped
- # reads unmapped
- # reads mapped to known contaminants
- # of uniquely mapped reads
- # duplicates
- expected read distribution pattern



Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

## Terminology





Graphic by Joern Toedling

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Read duplication

#### Origins for Read Duplication

- biological
- technical (e.g. PCR amplification, optical duplicates)



B Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

# Visualising duplication



Raw data quality control Mapping quality control Application specific quality control Batch effects

#### How to handle duplicate reads?

- DNA/ChIP-seq duplicate removal or estimation of biological duplication rate
- RNA-seq no duplication removal before analysis



Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

## Read distribution pattern: ChIP





Raw data quality control Mapping quality control Application specific quality control Batch effects

## Read distribution pattern: RNA-Seq





Mainz











Raw data quality control Mapping quality control Application specific quality control Batch effects

# Quality control of RNA-seq

- sequencing depth (unique mapping reads)
- rRNA content
- 5' to 3' distribution of reads (gene body coverage)
- strand specificity
- duplication rate
- distribution of reads over different gene classes



Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### RNA-seq: 5' to 3' coverage



Nastasja Kreim & Anke Busch - IMB Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

## RNA-seq: 5' to 3' coverage



Nastasja Kreim & Anke Busch – IMB Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

## RNA-seq: Contamination Screening





Raw data quality control Mapping quality control Application specific quality control Batch effects

#### RNA-seq: Counts on Features



Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### RNA-seq: duplication rate





graphic from Holger Klein based on DupRadar

Nastasja Kreim & Anke Busch – IMB Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

## RNA-seq: low complexity library





graphic from Holger Klein based on DupRadar

Nastasja Kreim & Anke Busch - IMB Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### RNA-seq: annotation issues

Se	tale 5 kb →	65.825.000		
	A STATE OF A	h de la company		N N O
	and the second sec	he stated at a state of the sta		<u>' / / / / / / / / / / / / / / / / / / /</u>
	And a state of the	An and the American sector		
	and the second se	Advance of the second states		-
	And the second	And the set the set of the set of the set		
	AND A REAL PROPERTY OF A DESCRIPTION OF A D	heading the account		
	and the second se			•
	and the second s	- Income and the second se	· · · · · · ·	-
	the state of the state of the state		1	-
	A REAL PROPERTY AND A REAL	An an anna an	• • •	-
	and the second state of th	a destado entre la calada de la	· · · · · · · · · · · · · · · · · · ·	-
	and the second sec	h de de la constante de la const	· · · ·	-
	and the second se	the state of the s		-
	and the state of t	A de terrer et et staat hants	• · · · · · · · ·	-
	THE REPORT OF THE PARTY OF THE	Ale division of the same his hold	• · · · •	-
	The Assessment of the Assessment of the Assessment of the	a des des des des services del serte		4
	AND ADDRESS AND ADDRESS OF ADDRESS AND ADDRESS	he had been a state of the state of the state		
110	THE REPORT OF THE PARTY OF THE	La Milda & Bhatta abluere		
	AND RECEIPTION OF THE DESIGN AND AND ALL ADDRESS.			•
	And the second se			-
2101	AND REAL PROPERTY AND ADDRESS OF A DESCRIPTION OF A DESCR	he with the section of the sector of the sec		ui .
$\gamma$	And the second sec	in the statistic designment in the second	a an tao a antar ar a	<b>M</b>
	and have a state of the state of the state of the state of the	ha strift i dine in string to		<u>u</u>
	and have been as an an and have a start of a start of the	ha fiting t. Bus de chtiete		
1702	256			



in house data (mm9), figure adapted from Emil Karaulanov

Nastasja Kreim & Anke Busch - IMB Mainz

Raw data quality control Mapping quality control Application specific quality control Batch effects

# RNA-seq: multi-mapping reads





in house data (mm9), figure adapted from Emil Karaulanov

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments







Raw data quality control Mapping quality control Application specific quality control Batch effects

## ChIP: input control



- biases by sonication
- large genomic variation (e.g. Aneuploidy, large InDels, CNV)
- artefacts of preparation



Rozowsky et al., Nature Biotech, 2009

Raw data quality control Mapping quality control Application specific quality control Batch effects

# ChIP-seq: enrichment quality control IPStrength



Fidel et al deepTools2: A next Generation Web Server for Deep-Sequencing Data Analysis. Nucleic Acids Molecular Biology Research (2016).

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

Raw data quality control Mapping quality control Application specific quality control Batch effects

# ChIP-seq: enrichment quality control IPStrength





INT

Fidel et al deepTools2: A next Generation Web Server for Deep-Sequencing Data Analysis. Nucleic Acids Molecular Biology Research (2016).

Nastasja Kreim & Anke Busch – IMB Mainz Design and Analysis of NGS Experiments

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Batch effects





Mainz

Source: http://www.molmine.com/magma/global\_analysis/batch\_effect.html

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Batch effects





Source: http://www.molmine.com/magma/global\_analysis/batch\_effect.html

Raw data quality control Mapping quality control Application specific quality control Batch effects

### Batch effects





Source: http://www.molmine.com/magma/global\_analysis/batch\_effect.html

Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Batch effects



Raw data quality control Mapping quality control Application specific quality control Batch effects

MB Mainz

#### Guidelines for Experiments

- Encode ChIP-seq, DNA-seq Guidelines
- Encode RNA-seq Standards
- ...



Encode Guidelines: https://www.encodeproject.org/data-standards/

Batch effects

### Summary II

• exploratory data analysis to ensure data quality

oinformatics



3 Mainz
### Summary II

- exploratory data analysis to ensure data quality
- use known unbiased quality control methods on the different analysis levels

Batch effects



## Summary II

- exploratory data analysis to ensure data quality
- use known unbiased quality control methods on the different analysis levels
- investigate the similarities between samples (principal component analysis and clustering)

Batch effects



## Summary II

- exploratory data analysis to ensure data quality
- use known unbiased quality control methods on the different analysis levels
- investigate the similarities between samples (principal component analysis and clustering)

Batch effects

• visualise your data



Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Tools

- Rawdata: FastQC
- Mapper: Bowtie, BWA, STAR, HiSat2
- Mapping-QC: samtools, Picard tools, qualimap, deepTools
- Duplication: bamUtils, DupRadar
- RNA diff. expression: DESeq2, edgeR
- RNA-seq QC: RNASeqQc, RNAQC
- ChIP-seq QC: deepTools
- Peak calling: MACS2
- Contamination: BLAST, FastqScreen
- Misc Analysis: SeqMonk
- Visualisation: IGV, SeqMonk, UCSC genome browser, Washington epigenome browser



Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Acknowledgements

- Holger Klein
- Bioinformatics Core Facility:
  - Emil Karaulanov
  - Fridolin Kielisch
  - Martin Oti
  - Giuseppe Petrosino
  - Frank Rühle
  - Sergi Sayols Puig



Raw data quality control Mapping quality control Application specific quality control Batch effects

#### Acknowledgements

- Holger Klein
- Bioinformatics Core Facility:
  - Emil Karaulanov
  - Fridolin Kielisch
  - Martin Oti
  - Giuseppe Petrosino
  - Frank Rühle
  - Sergi Sayols Puig

# Thank you for your attention.

