

Fusion V5 (Gene quantification)

Bioinformatics Development

Sandy

2025.03.24

QC metrics overview

Ref. issue:

- <https://actg.atlassian.net/browse/ABIE-971>

413 target exons:

- /mnt/RD_Develop/sandyteng/ACTFusionV5/code/fusionv4_annoloci2bed_test/targetexonbed/fusionv4.MANE.v0.95.GENCODE.r38.candidate.exons.transcript.bed

- Tools & fusion workflows

	STAR (arriba's workflow: STAR + arriba)	Fusion v4 (bwa-based)
Alignment analysis	STAR (to genome)	bwa-mem (to preferred transcriptome, MANE, GENCODE-r38)
(I) # of primary mapped reads	samtools flagstats (~81.7% from Twist NextSeq data)	samtools flagstats (~88.6% from Twist NextSeq data)
(II) % of on-target/probe-anchored reads	calculate_probe_reads.sh (in-house utility: samtools + bedtools) (~54% On-Target reads, Twist NextSeq data)	calculate_probe_reads.sh (in-house utility: samtools + bedtools) To-do
Read trimming	NA	trimadap
Counting (expression quantification)	HTseq ("htseq-count"), FeatureCounts ("featureCounts")	quantify_preferred_exons.v2.py 1. (transcript-level) via "htseq-count" => Need gtf file for preferred transcripts => Some arguments are not applicable for bwa (no 'NH' tag) 2. (transcript-level) obtain alignments from *callingresult.txt file for each sample => Use "WILDTYPE" reads produced by the caller to quantify gene expression
Fragmentation size	NA	fastp (insert size → peak, source file: *.fastp.merge.json) (129-153 bp insertion size, Twist NextSeq data)
Duplication rate	NA	fastp (duplication → rate, source file: *.fastp.merge.json) (29%-34% duplication rate, Twist NextSeq data)

Counting (expression quantification)

- Quantification scenarios
 - Htseq (+ arriba.STAR.bam)
 - FeatureCounts (+ arriba.STAR.bam)
 - quantify_preferred_exons.v2.py (in-house script) (+ fusionv4.bwa.bam)
- Analysis workflow
 - Gene count quantification (via htseq, featurecounts, quantify_preferred_exons.v2.py)
 - Target gene count extraction (only compare the 220 target genes defined in twist.covered.bed (via grep -wf))
- Result summary
 - Gene count obtained from htseq and featurecounts are similar (correlation 99.9%)
 - Gene count quantified from fusion v4 and arriba workflows are similar (correlation 99.3%)

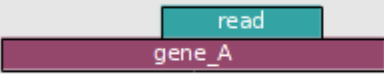
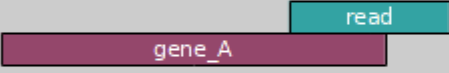


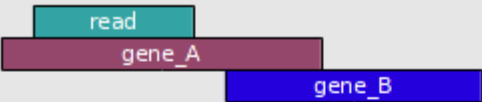
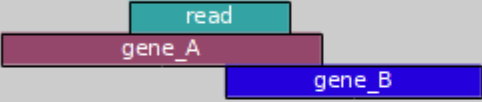
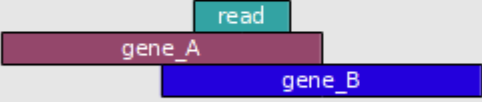
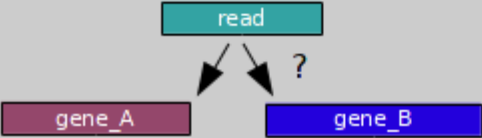
Target gene:

- /mnt/RD_Develop/sandyteng/workdir/bed_intersect/Twist_kit/Covered_Regions_RNA_Fusions_4X_TE-98493102_GRCh38.gene.list.txt (Twist)

HTSeq-count

- Default options for feature count (gene count)
- **-t exon**
(default feature type => 3rd column in GTF file)
- **-i gene_id**
(default id attribute => feature ID)
- **-m union**
(default read overlapping handling)
- **--nonunique none**
(default mode for reads aligned to more than one feature in the “-m” option)

Ref. link
 • [Htseq-count docs](#)

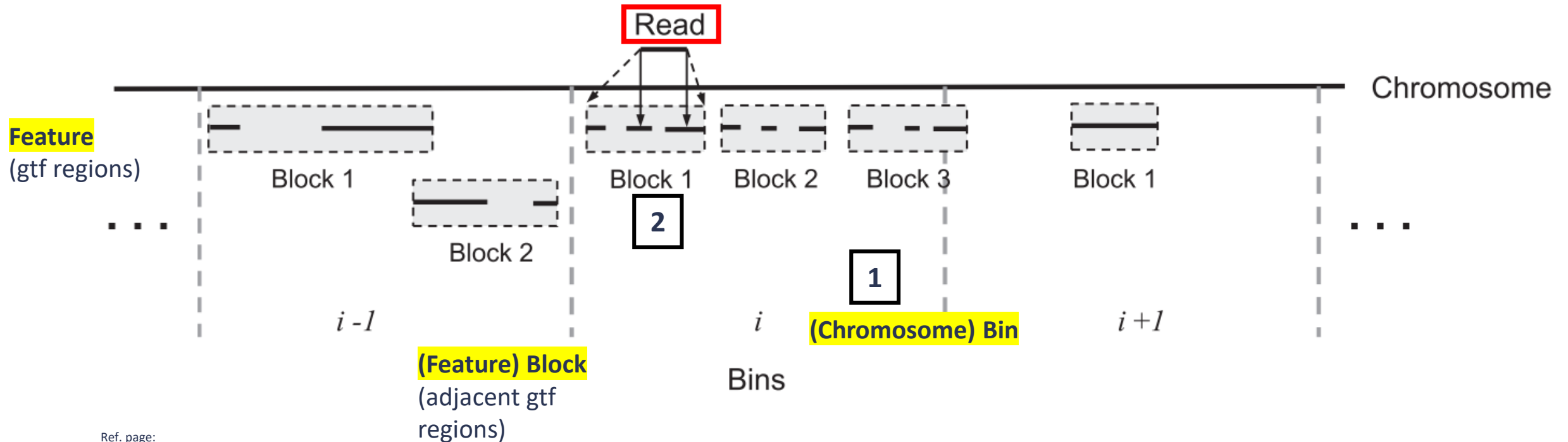
	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous (both genes with --nonunique all)	gene_A	gene_A
	ambiguous (both genes with --nonunique all)		
	alignment_not_unique (both genes with --nonunique all)		

FeatureCounts

- FeatureCounts (algorithm)

Steps:

- Overlap of reads with features
- Multiple overlaps
- Chromosome hashing
- Genome bins and feature blocks



Ref. page:

- [featureCounts: an efficient general purpose program for assigning sequence reads to genomic features](#)

Exon number vs Count

- MET

gene count

gene_id count

MET **18002** => identical to # of WILDTYPE MET reads

exon count

=>

(-i) *callingresult.txt ("ENST00000397752.8" => "MET")

167 WILDTYPE MET:15,16

280 WILDTYPE MET:15,16,17

53 WILDTYPE MET:16

372 WILDTYPE MET:16,17

12 WILDTYPE MET:16,17,18

=> (-o) exon.count ("ENST00000397752.8" exon 16)

=> ENST00000397752.8 16 884 (167+280+53+372+12 = 884)

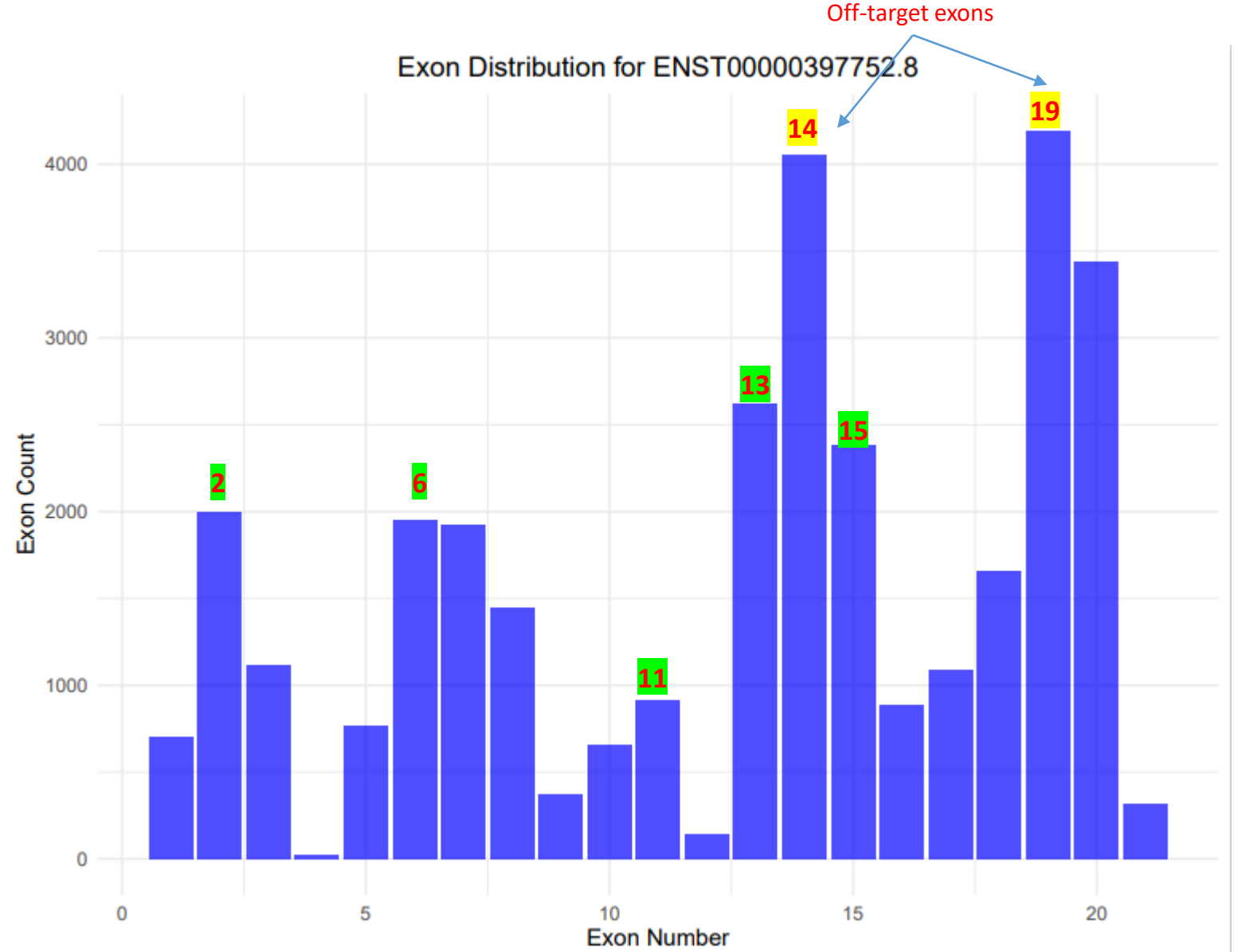


Exon number vs Count

• MET

chr7	116699071	116700284	MET-exon-fusionv4-2
chr7	116755355	116755515	MET-exon-fusionv4-6
chr7	116763050	116763268	MET-exon-fusionv4-11
chr7	116771498	116771654	MET-exon-fusionv4-13
chr7	116774881	116775111	MET-exon-fusionv4-15

gene_id count	ENST00000397752.8	1	707
MET 18002	ENST00000397752.8	2	2000
=> identical to	ENST00000397752.8	3	1117
# of WILDTYPE MET reads	ENST00000397752.8	4	27
	ENST00000397752.8	5	763
	ENST00000397752.8	6	1952
	ENST00000397752.8	7	1928
	ENST00000397752.8	8	1443
	ENST00000397752.8	9	373
	ENST00000397752.8	10	653
	ENST00000397752.8	11	918
	ENST00000397752.8	12	145
	ENST00000397752.8	13	2624
	ENST00000397752.8	14	4055
	ENST00000397752.8	15	2385
	ENST00000397752.8	16	884
	ENST00000397752.8	17	1088
	ENST00000397752.8	18	1655
	ENST00000397752.8	19	4196
	ENST00000397752.8	20	3443
	ENST00000397752.8	21	314



Gene quantification (FusionV4)

Ref. issue

<https://actg.atlassian.net/browse/ABIE-987>

<https://actg.atlassian.net/browse/ABIE-988>

- FusionV4 processes
 - R1.fq.gz, R2.fq.gz (input files) -> mergefastq ("mergefastq") -> trimadap -> fastp -> bwaisoform -> **bwase -> fusioncalling** -> fuscalle2QC
- **bwase ("bwa") -> fusioncalling ("ACTGfuscall.py") -> quantifygene ("quantify_preferred_exons.v2.py")**
 - bwase
 - Input files: preferred.transcriptome.fasta, preferred.transcriptome.fasta.indices
 - Output files: aligned.fusionv4.bam, aligned.fusionv4.bam.bai
 - Fusioncalling
 - Input files:
aligned.fusionv4.bam, preferred.transcriptome.exons.annotation,
protein.fasta, protein.fasta.meta, qc.thresholds.config
 - Output files:
callingresult.txt, gspcallingresult.txt, protein_seq.meta.txt, callingform.txt
 - quantifygene
 - Input file: callingresult.txt
 - Output files: gene.count, exon.count

Gene quantification (Arriba)

Ref. issue

<https://actg.atlassian.net/browse/ABIE-987>

<https://actg.atlassian.net/browse/ABIE-988>

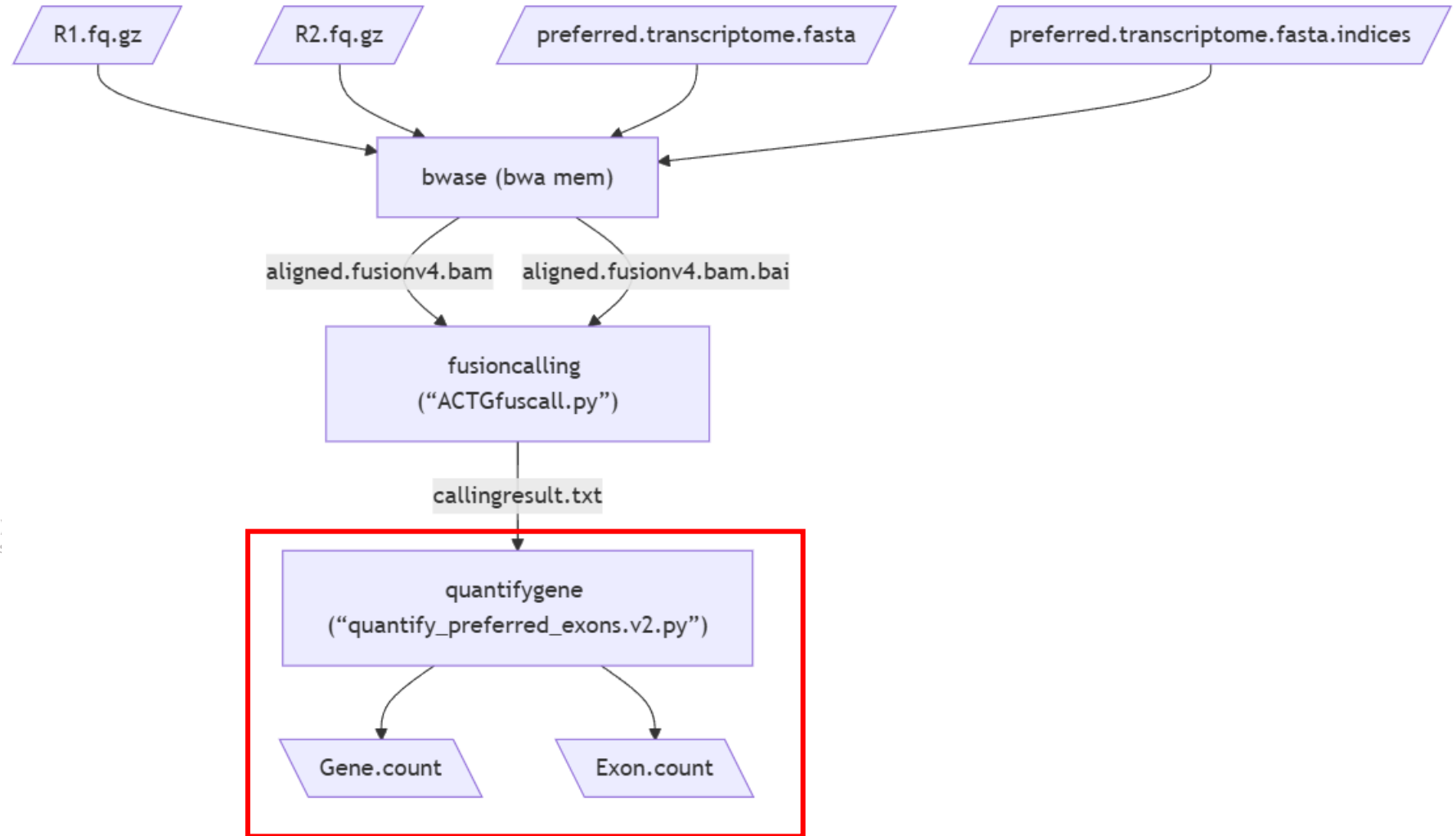
- Arriba processes
 - R1.fq.gz, R2.fq.gz (input files) -> **STAR** -> Arriba
- **STAR ("STAR") -> quantifygene ("HTSeq" / "featureCounts")**
 - STAR
 - Input files: GRCh38.fa, GRCh38.fa.indices, RefSeq_hg38.gtf
 - Output files: aligned.arriba.bam
 - quantifygene (htseq)
 - Input files: aligned.arriba.bam, RefSeq_hg38.gtf
 - Output files: genes_htseq.count (gene.count)
 - quantifygene (featureCounts)
 - Input files: aligned.arriba.bam, RefSeq_hg38.gtf
 - Output files: genes_featureCounts.count (gene.count), genes_featureCounts.count.summary (gene.count.summary)

Workflow

- FusionV4

```
graph TD;
%% Initial Inputs
I1[/R1.fq.gz/]
I2[/R2.fq.gz/]
I7[/preferred.transcriptome.fasta/]
I8[/preferred.transcriptome.fasta.indices/]
%% FusionV4 Workflow
I1 --> A1
I2 --> A1
I7 --> A1["bwase (bwa mem)"]
I8 --> A1["bwase (bwa mem)"]
A1 --> B1
A1 --> B1["fusioncalling (\"ACTGfuscall.py\")"]
B1 --> C1["callingresult.txt"]
C1 --> D1
D1 --> E1[/Gene.count/]
D1 --> E2[/Exon.count/]
style D1 stroke:#f00,stroke-width:2px
```

graph TD;
%% Initial Inputs
I1[/R1.fq.gz/]
I2[/R2.fq.gz/]
I7[/preferred.transcriptome.fasta/]
I8[/preferred.transcriptome.fasta.indices/]
%% FusionV4 Workflow
I1 --> A1
I2 --> A1
I7 --> A1["bwase (bwa mem)"]
I8 --> A1["bwase (bwa mem)"]
A1 --> B1
A1 --> B1["fusioncalling (\"ACTGfuscall.py\")"]
B1 --> C1["callingresult.txt"]
C1 --> D1
D1 --> E1[/Gene.count/]
D1 --> E2[/Exon.count/]
style D1 stroke:#f00,stroke-width:2px

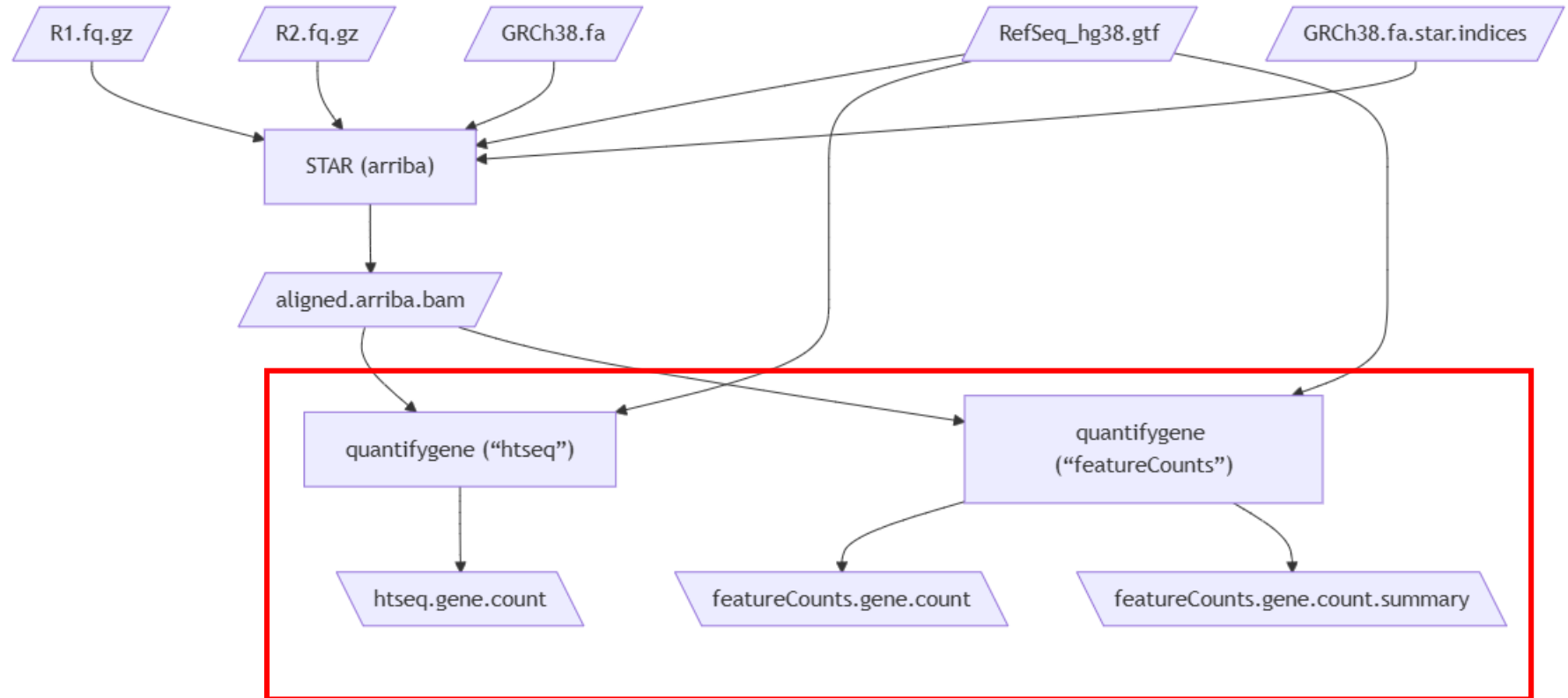


Workflow

- Arriba

```
graph TD;
  %% Initial Inputs
  I1[/R1.fq.gz/]
  I2[/R2.fq.gz/]
  I6[/GRCh38.fa/]
  I7[/GRCh38.fa.star.indices/]
  I8[/RefSeq_hg38.gtf/]
  %% FusionV4 Workflow
  I1 --> A1
  I2 --> A1
  I6 --> A1
  I7 --> A1["STAR (arriba)"]
  I8 --> A1
  A1 --> I9[/aligned.arriba.bam/]
  I9 --> B1["quantifygene ('htseq')"]
  I9 --> C1["quantifygene ('featureCounts')"]
  %%A1 --> |aligned.arriba.bam| B1["quantifygene ('htseq')"]
  B1 --> O1[/htseq.gene.count/]
  %%A1 --> |aligned.arriba.bam| C1["quantifygene ('featureCounts')"]
  C1 --> O2[/featureCounts.gene.count/]
  C1 --> O3[/featureCounts.gene.count.summary/]
  style B1 fill:#d9d9ff,stroke:#333,stroke-width:1px
  style C1 fill:#d9d9ff,stroke:#333,stroke-width:1px
  style O1 fill:#d9d9ff,stroke:#333,stroke-width:1px
  style O2 fill:#d9d9ff,stroke:#333,stroke-width:1px
  style O3 fill:#d9d9ff,stroke:#333,stroke-width:1px
```

graph TD;
 %% Initial Inputs
 I1[/R1.fq.gz/]
 I2[/R2.fq.gz/]
 I6[/GRCh38.fa/]
 I7[/GRCh38.fa.star.indices/]
 I8[/RefSeq_hg38.gtf/]
 %% FusionV4 Workflow
 I1 --> A1
 I2 --> A1
 I6 --> A1
 I7 --> A1["STAR (arriba)"]
 I8 --> A1
 A1 --> I9[/aligned.arriba.bam/]
 I9 --> B1["quantifygene ('htseq')"]
 I9 --> C1["quantifygene ('featureCounts')"]
 %%A1 --> |aligned.arriba.bam| B1["quantifygene ('htseq')"]
 B1 --> O1[/htseq.gene.count/]
 %%A1 --> |aligned.arriba.bam| C1["quantifygene ('featureCounts')"]
 C1 --> O2[/featureCounts.gene.count/]
 C1 --> O3[/featureCounts.gene.count.summary/]

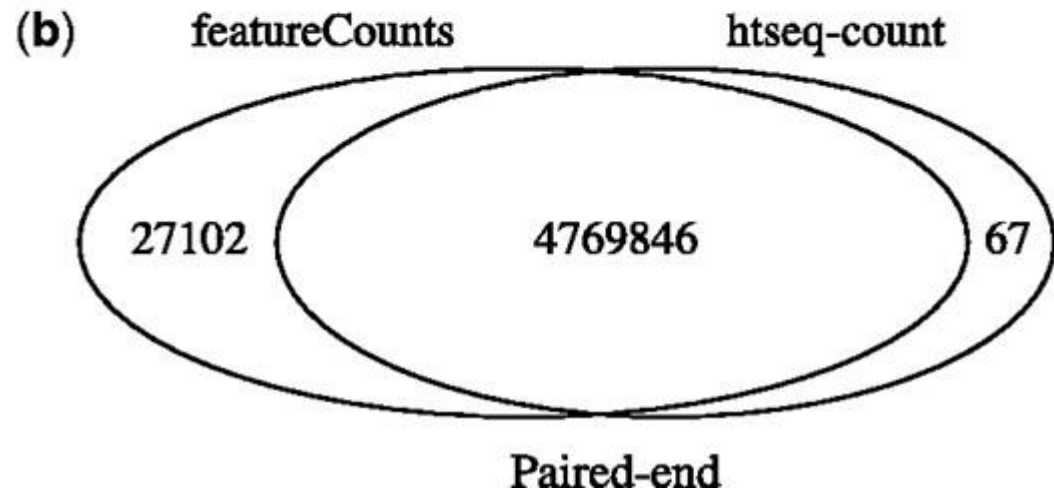
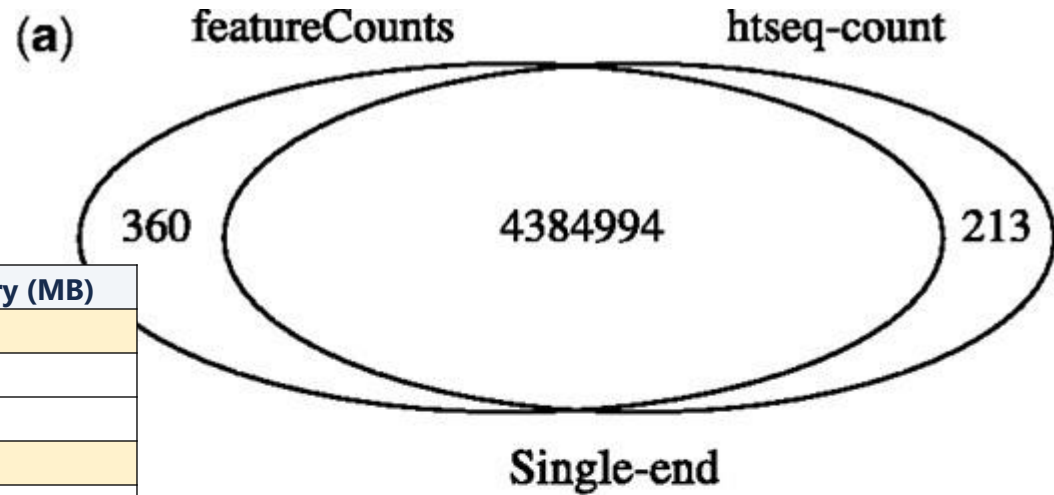


Gene count comparison (paper)

- HTSeq/FeatureCounts

Method	Number of fragments	Time (min)	Memory (MB)
<i>featureCounts</i>	5 392 155	0.9	4
<i>CountOverlaps (whole genome at once)</i>	5 392 155	24.4	7000
<i>CountOverlaps (by chromosome)</i>	5 392 155	36.6	783
<i>htseq-count (union)</i>	4 978 050	36	31
<i>htseq-count (intersection-nonempty)</i>	4 993 644	35.7	31
<i>coverageBED</i>	5 366 902	4.4	41

Ref:
[featureCounts: an efficient general purpose program for assigning sequence reads to genomic features](#)



Gene count comparison (AANB02_202_AD02_AA-23-08153)

- HTSeq/FeatureCounts
- FeatureCounts/quantify_preferred_exons.v2.py

Pearson's r (B, C) (featureCounts, htseq)	0.999996602
n	244
t-value $(r \cdot \sqrt{n-2}) / (\sqrt{1-r^2})$	5967.755077
p-value TDIST(x, deg_freedom, tails)	0
Pearson's r (C, D) (featureCounts, quantify_preferred_exons.v2.py)	0.993358852
n	244
t-value	134.307269
p-value	2.8426E-229

Ref. issue
<https://actg.atlassian.net/browse/ABIE-987>
<https://actg.atlassian.net/browse/ABIE-988>

- Ref:
- HTSeq/FeatureCounts
 - bedtools coverage/samtools depth
 - quantify_preferred_exons.v2.py

Tool overview

	HTSeq/FeatureCounts	bedtools coverage /samtools depth	quantify_preferred_exons.v3.py
Quantification level	Gene Level (predefined intervals within gtf => gene id recognition)	Base Level (bedtools coverage -d /samtools depth) Interval Level (bedtools coverage)	Gene Level + Exon Level
Limitations	<p>Some arguments are not applicable for bwa (no "NH" tag)</p> <p>Count gene using the predefined gtf (merged the same gene_id) => limit to predefined gene intervals (may not encompass MANE 1.4 transcripts)</p>	<p>samtools depth is preferred for SAM FLAG sensitivity (duplication removal)</p> <p>ABIE-976: "bedtools coverage" vs. "samtools depth"</p> <p>Done</p> <p>bedtools coverage fails to identify read fragment => extra care is required for result interpretation see details for https://github.com/ACTGenomics/panel_gene_coverage</p>	<p>Only quantify exons defined in the preferred transcripts (MANE 0.95 + GENCODE-r38) => one may change the preferred transcripts to MANE 1.4</p> <p>Rely on fusion v4 calling result => only work for fusion v4 pipeline</p>

Ref. issue
<https://actg.atlassian.net/browse/ABIE-987>
<https://actg.atlassian.net/browse/ABIE-988>

Make
Personalized
Medicine
Accessible to All

