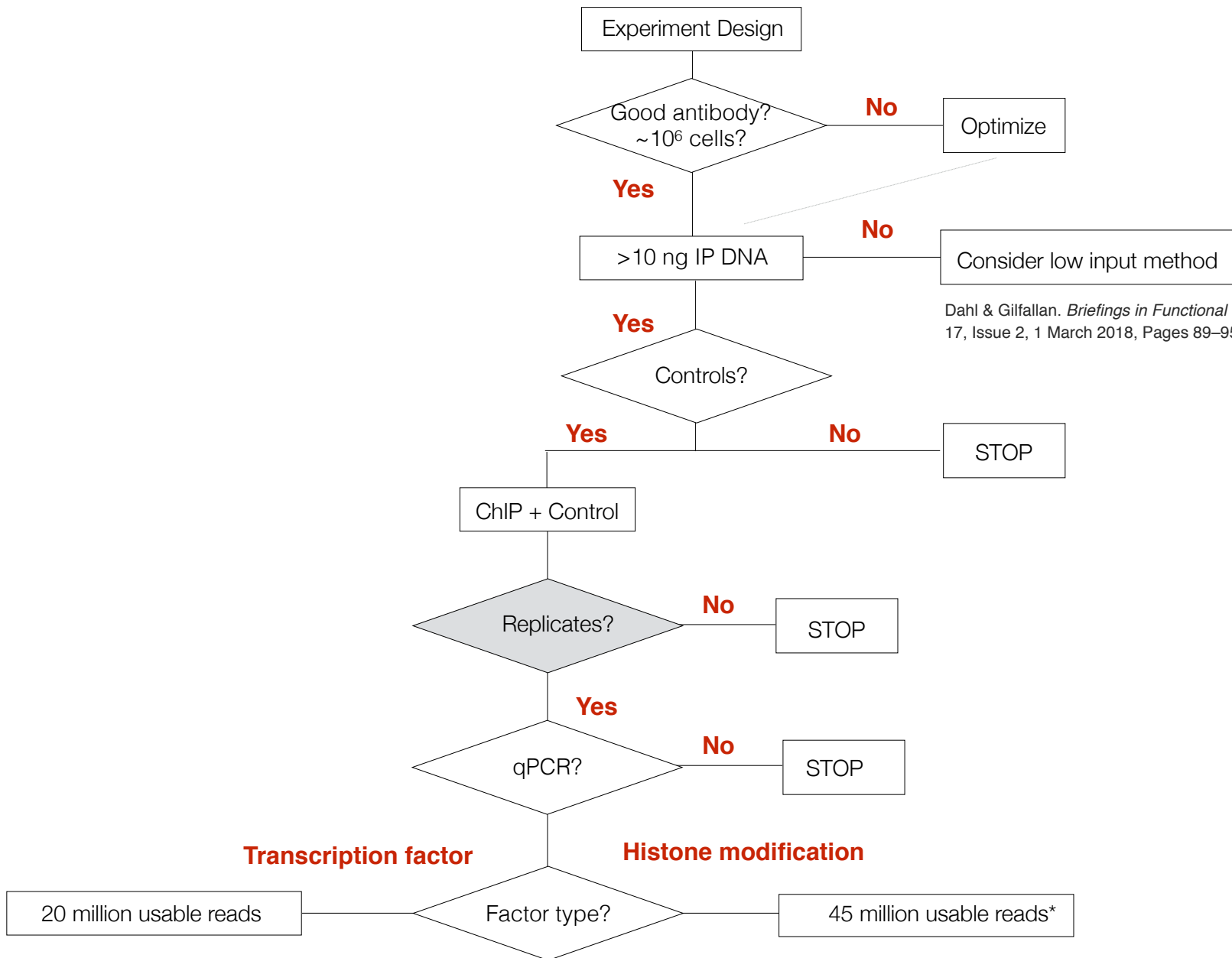


ChIP-seq Analysis Workflow and Troubleshooting

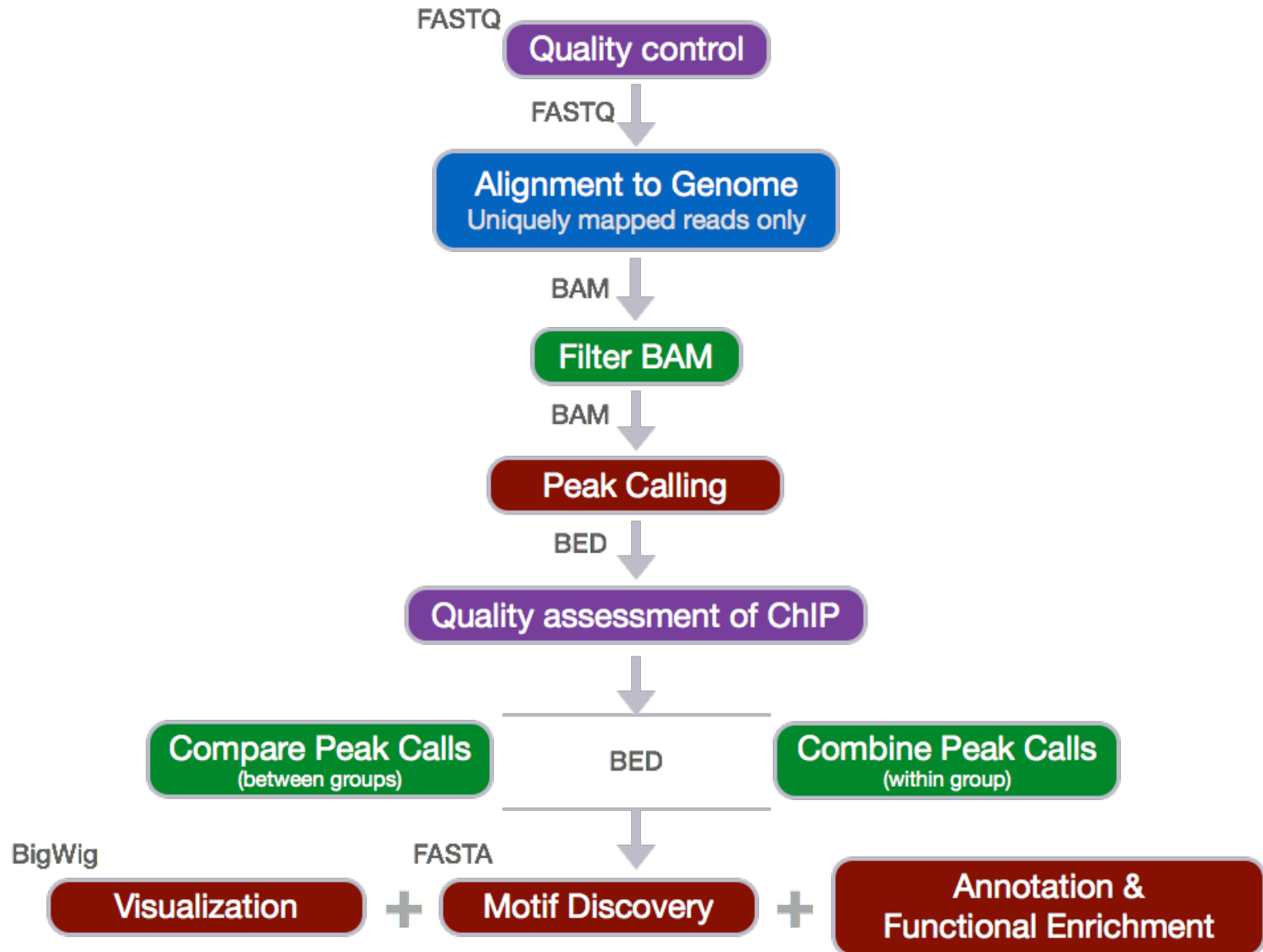
Before the sequencer...

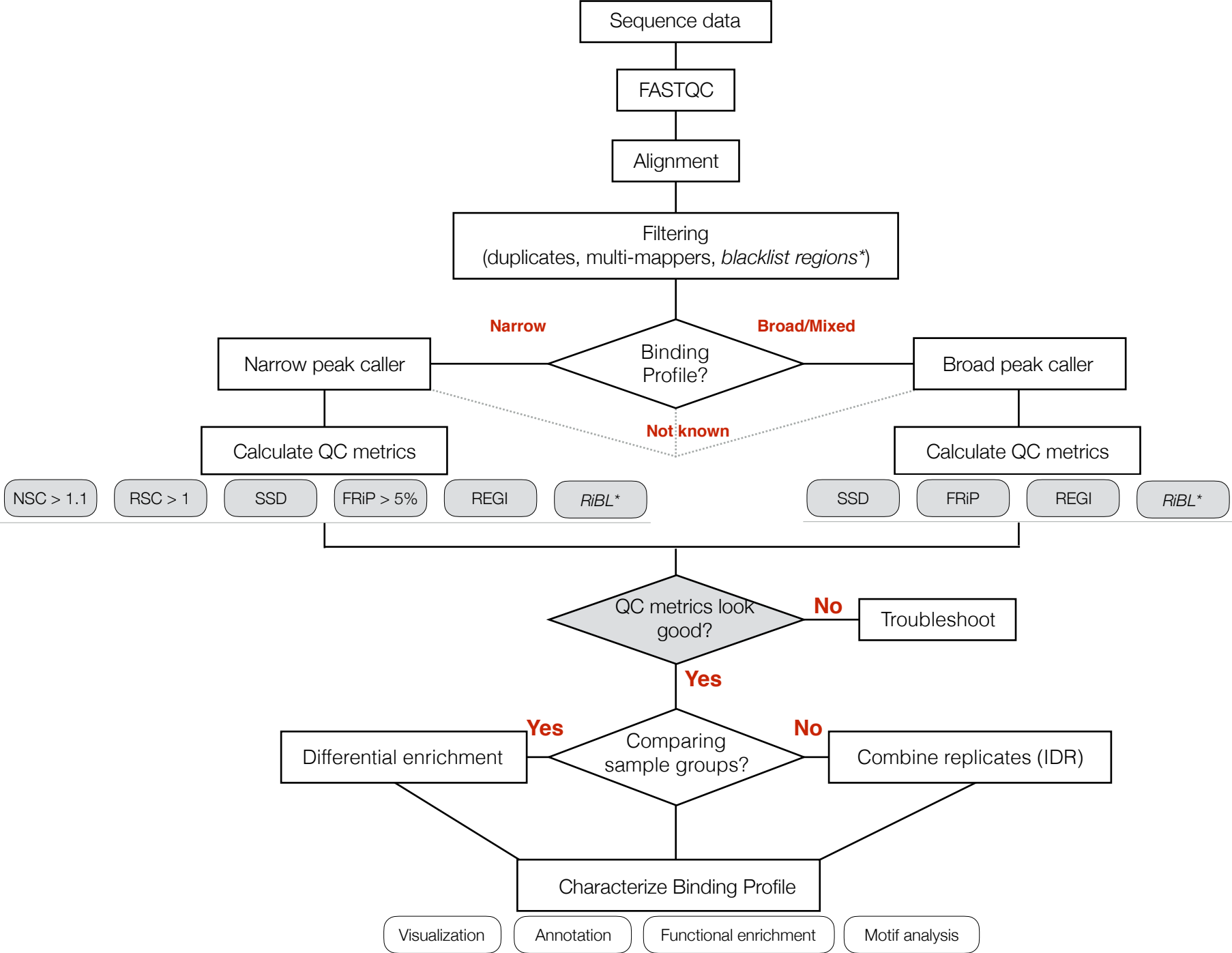


Dahl & Gilfillan. *Briefings in Functional Genomics*, Volume 17, Issue 2, 1 March 2018, Pages 89–95,

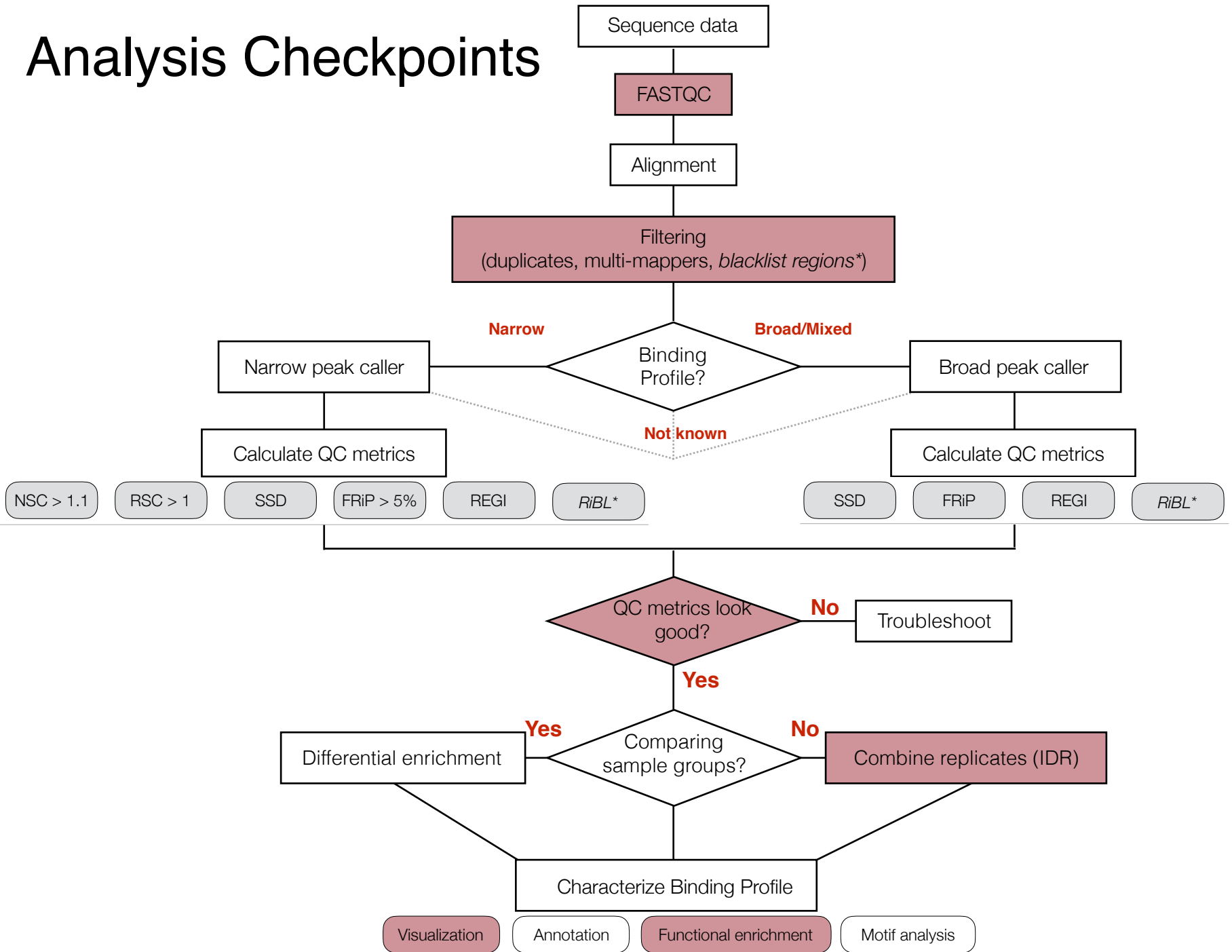
After sequencing...

ChIP-seq Workflow



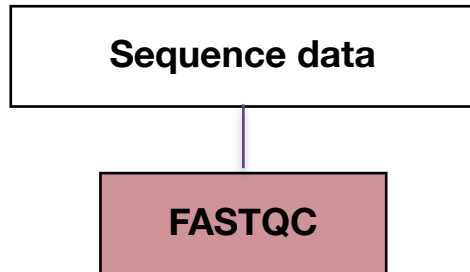


Analysis Checkpoints



Quality Checks: Raw Data

The quality checks at this stage in the workflow include:



1. Checking the **quality of the base calls** to ensure that there were no issues during sequencing
2. Examining the reads to ensure their **quality metrics adhere to our expectations** for our experiment
3. Exploring reads for **contamination**

✔ **Per base sequence quality**

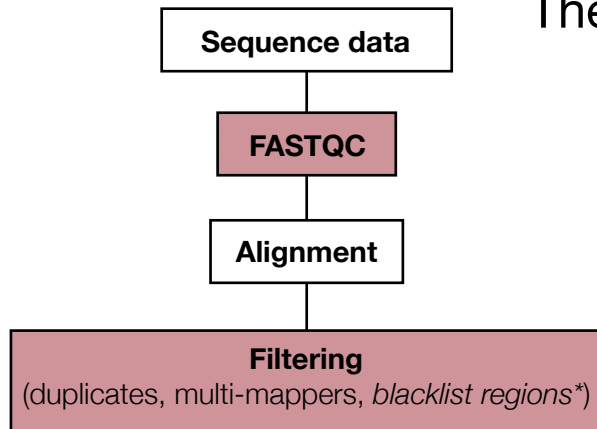
[illegible]

Quality Checks: Raw Data

Troubleshooting raw data quality problems:

- Low sequence quality reads
 - loss of signal in later cycles, technical problems with sequencer
- Unexpected %GC for organism
 - contaminating sequences: different species, adapters, vector
- High level of sequence duplications
 - low complexity library, too many cycles of PCR amplification / too little starting material, biological duplicates
- Over-represented sequences can be biologically significant or represent bias
 - sequences that represent binding sites
 - contaminating sequences: adapters, vector

Quality Checks: Aligned Data



The quality checks at this stage in the workflow include:

1. Checking the **total percent of reads aligning** to the genome
2. Determine the **percent of duplicate reads**
3. Determining the **percent uniquely mapping** reads
4. Identify percent of reads **mapping in blacklist regions**
5. Checking percent of paired-end reads that are **properly paired**

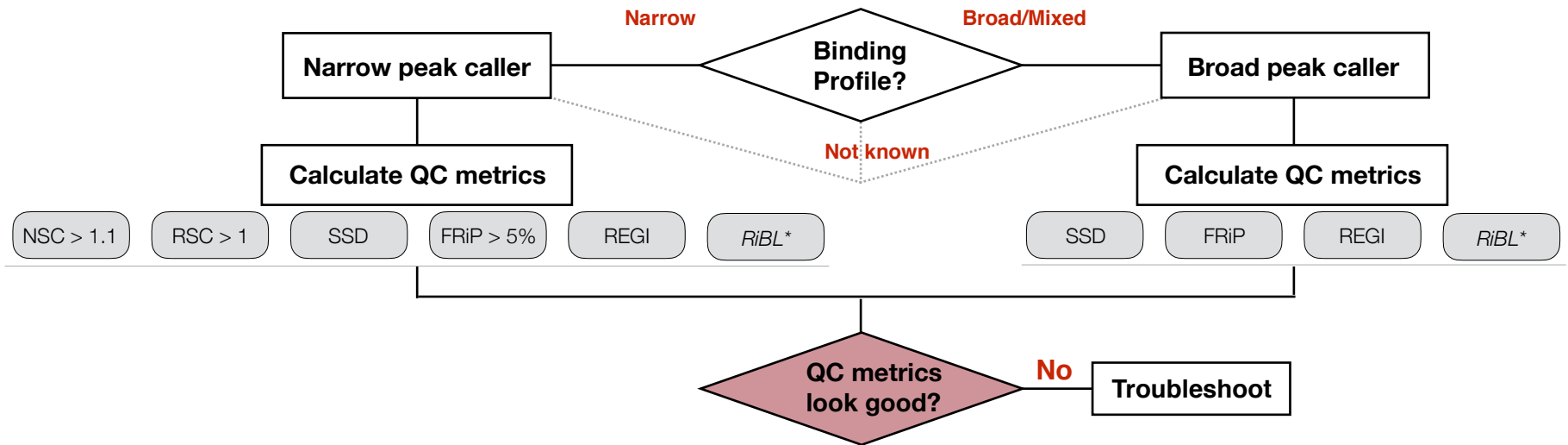
Quality Checks: Aligned Data

Troubleshooting aligned data quality problems:

- Low percentage (< 70%) of **total reads** aligned
 - poor quality reads, contaminating sequences, inappropriate alignment parameters chosen, poor quality reference genome
- Low percentage (< 60%) of **uniquely** aligning reads
 - high number of multi-mappers, high number of duplicates (20-30% duplication rate)
- High percentage (> 10%) of reads mapping in **blacklist regions**
- For paired-end data: large number of reads **not properly paired**
 - poor quality reads

*Even if your samples meet the suggested thresholds, **always filter** for (duplicates), multi-mappers and blacklist regions.*

Quality Checks: Peak Calling



Quality Checks: Peak Calling

The quality checks at this stage in the workflow include:

1. Evaluate **degree of enrichment** (FRiP, coverage plots, SSD)
2. Evaluating **signal-to-noise** using strand cross-correlation based metrics (NSC, RSC, fragment length)
3. Evaluate enrichment **within specific genomic regions/features**, and within **known artifact regions**
4. **Compare and contrast** measures described above based on **thresholds** and/or **what you anticipate** for the binding profile

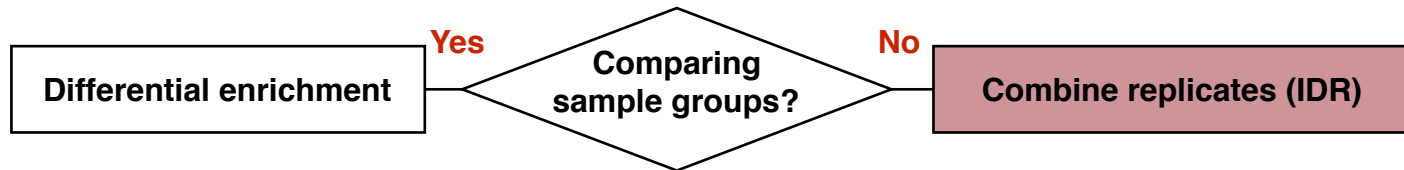
Quality Checks: Peak Calling

Troubleshooting ChIP quality problems:

- Low FRiP ($< 5\%$), Low SSD scores
 - poor enrichment due to poor antibody, the protein of interest binds few regions
- Low **signal-to-noise** ($RSC < 0.8$, $NSC < 1.1$), bad fragment length estimates, low diversity of depth
 - IP did not work, high cross-reactivity of antibody with other proteins
- No enrichment in anticipated **genomic features**
 - experiment did not work, other regions identify interesting novel behaviors
- High RiBL ($> 10\%$)
 - repeat regions or other artifact regions displaying artificial enrichment, this can drive up the SSD score

Evaluate these metrics collectively for each sample.

Quality Checks: Handling Replicates



The quality checks at this stage in the workflow include:

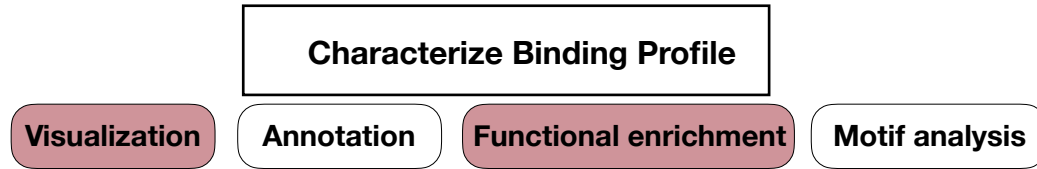
1. Using bedtools to perform a crude **comparison of peak calls** between replicates
2. **Statistical evaluation** of reproducibility between replicates using the IDR pipeline

Quality Checks: Handling Replicates

Replicate Evaluation Goals:

- Ensure that there is concordance in the peaks being called across all replicates in a sample group
- For IDR analysis use a more liberal threshold for peak calling to increase the search space
- Focus on peaks that meet the $IDR < 0.05$ threshold for downstream analysis (visualization and functional analysis)
- Pseudo-replicate analysis for reporting in a publication

Quality Checks: Characterizing the binding profile



The quality checks at this stage in the workflow include:

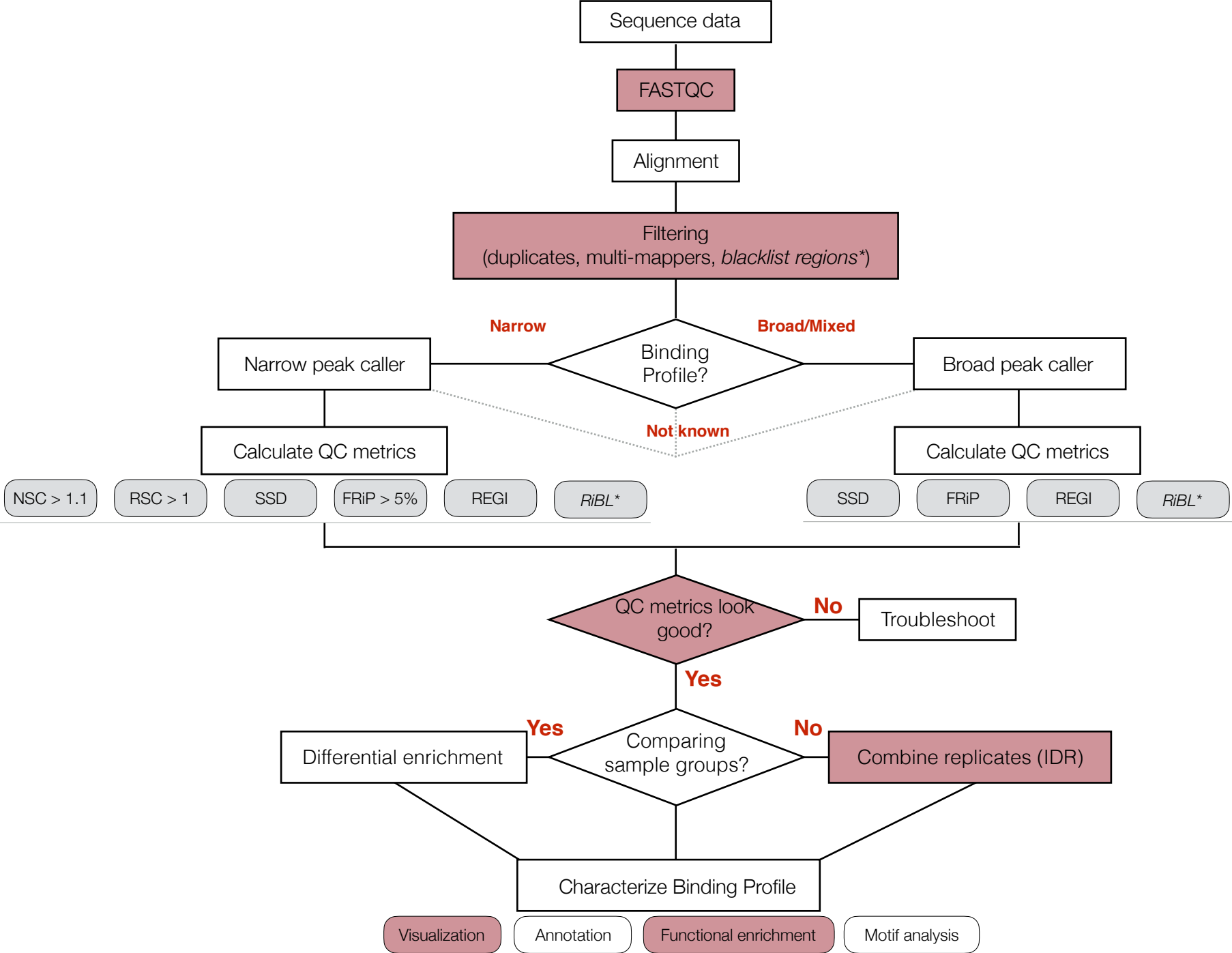
1. Using deepTools to **visualize the data** and evaluate enrichment in specific regions.
2. Look at **specific regions of interest** individually (anticipated target genes). Overlay relevant public datasets (using IGV or a genome viewer).
3. Evaluate target genes and assess **functional analysis** to see if there is any biological relevance.

Quality Checks:

Characterizing the binding profile

Troubleshooting problems:

- Do we see enrichment associated with genomic features we anticipated?
- Does the visual inspection validate what we know/identified from statistical analysis?
- Do these target genes collectively represent specific pathway(s)? Is there significant over-representation of certain biological processes?
- Is this all relevant based on what I know about my protein of interest?



ChIP-seq Workflow

