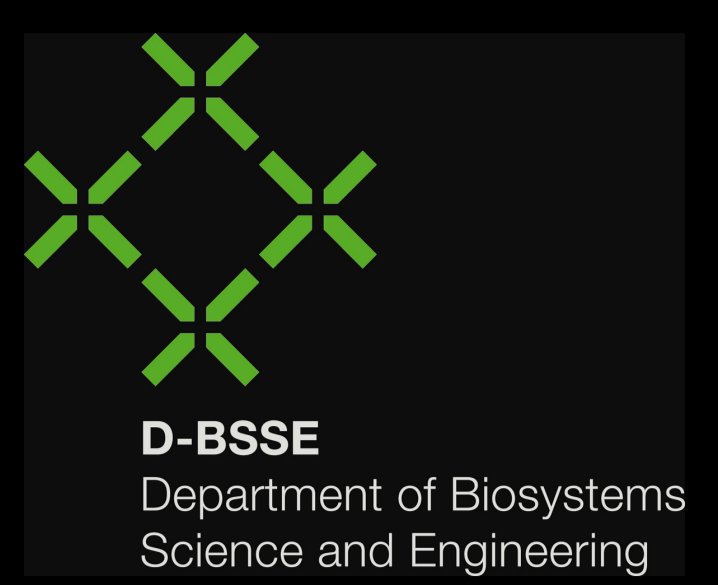


Deep Sequencing

Werner Van Belle¹, Ina Nissen¹, Michael Stadler² and Christian Beisel¹



(1) Laboratory for Quantitative Genomics, Department Biosystems Science and Engineering - ETH Zurich
(2) Friedrich Miescher Institute for Biomedical Research

Abstract

In parallel to the human genome sequencing initiative several new technologies have emerged that allow sequencing at unprecedented throughput and low costs. These approaches are generally referred to as "deep sequencing". They enable researchers to not only re-sequence genomes and thus to identify genome variations but also to quantify the abundance of experimentally enriched fractions of the genome. The very large numbers of short individual sequence reads produced by the Illumina Genome Analyzer (currently approx. 50 million reads per instrument run) are well suited to make direct quantitative measurements of the sequence content of a DNA sample. By determining a short sequence read from each of many randomly selected molecules from the sample and then mapping each sequence read onto the reference genome, the identity of each starting molecule is learned, and its frequency in the sample can be calculated. Desired levels of sensitivity and statistical certainty, needed to detect rare molecular species, can be achieved by adjusting the total number of sequence reads. Sequence census assays do not require knowing in advance that a sequence is of interest as a promoter, enhancer or RNA-coding domain, as most current microarray designs do. The combination of chromatin immunoprecipitation assays with the subsequent quantitative analysis of the enriched DNA sample by deep sequencing (ChIP-seq) has been proven to be of great value for whole mammalian genome approaches in several high-profile studies published over the last year. At D-B SSE we have established a deep sequencing unit based on Illumina sequencing technology located in the new "Laboratory for Quantitative Genomics". This poster gives an overview of the sequencing technology and the data analysis pipeline. Furthermore it provides insights into the quality of our functional genomics data recently generated by ChIP-seq and RNA-seq.

Workflow and Applications

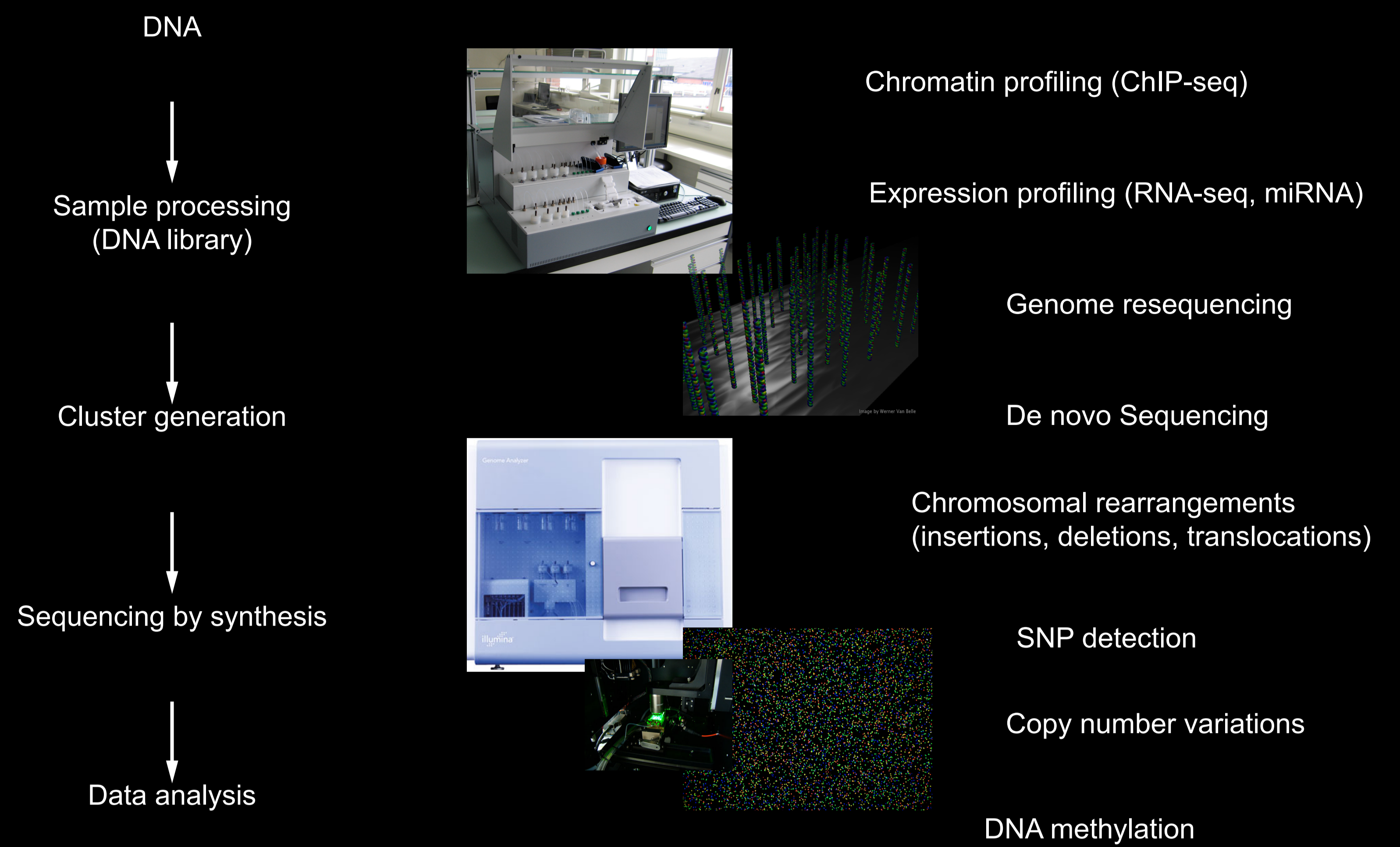


Image Acquisition

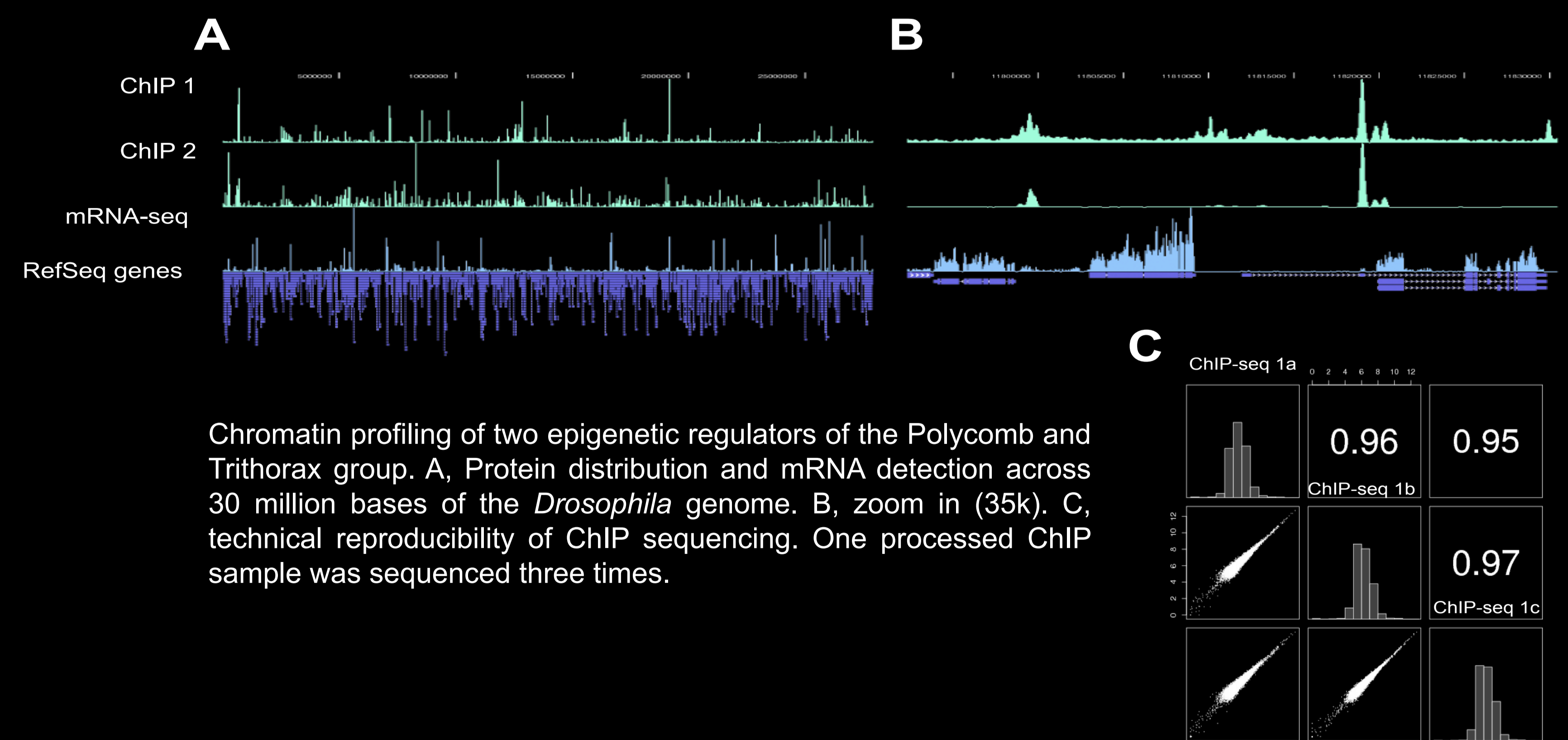


Short Fragment Alignment



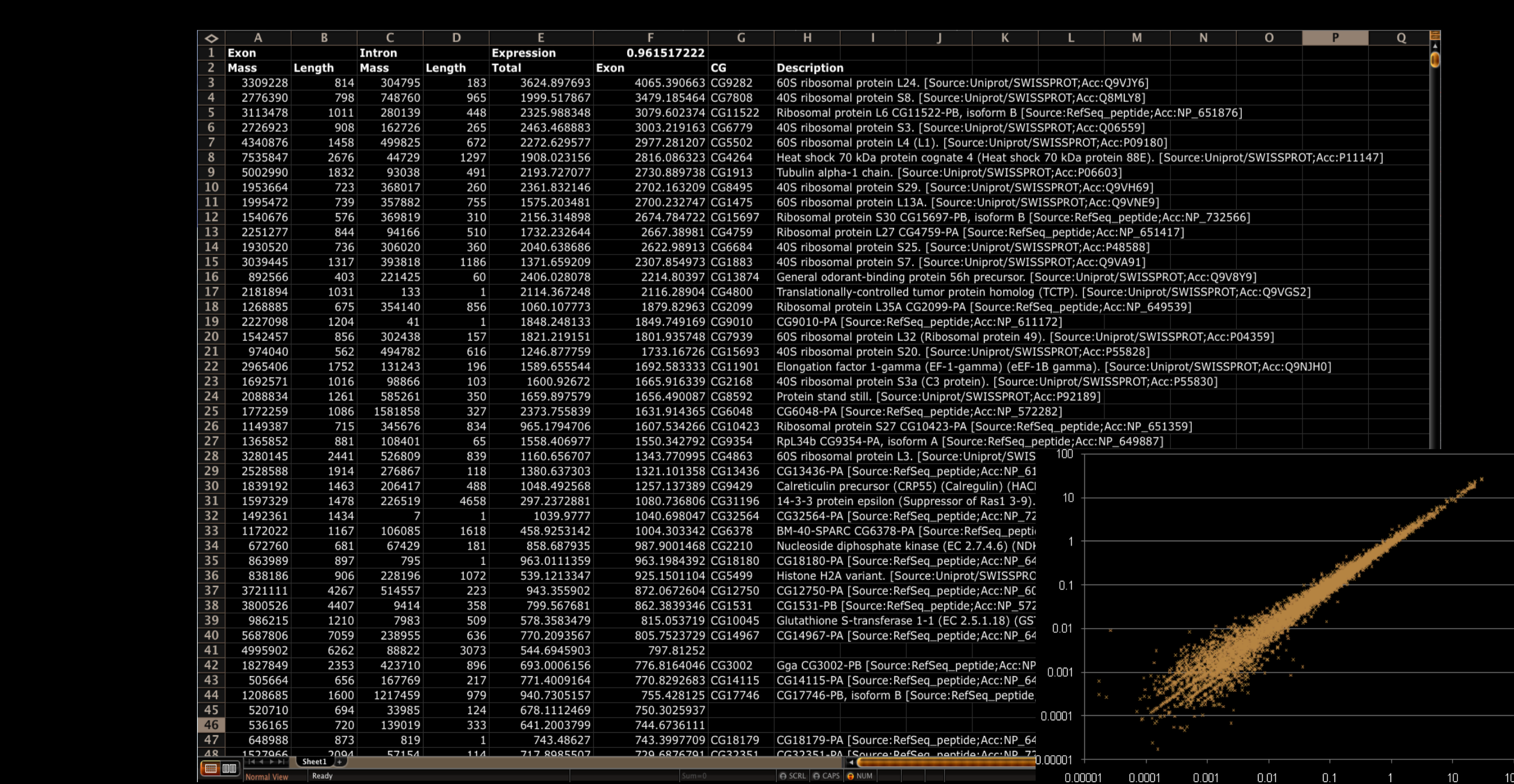
The sequenced short fragments can be aligned to a reference genome (yellow). The alignment program "Eland" allows for at most two mismatches per fragment. "PhageAlign" finds the best match for each fragment.

Chromatin Profiling



Chromatin profiling of two epigenetic regulators of the Polycomb and Trithorax group. A, Protein distribution and mRNA detection across 30 million bases of the *Drosophila* genome. B, zoom in (35k). C, technical reproducibility of ChIP sequencing. One processed ChIP sample was sequenced three times.

Expression Profiling



Data Delivery

Images: 100 Gb per lane - 800 Gb per flowcell
IPAR Output: 10.4 Gb/lane - 83.2 Gb/flowcell
Intensity files: 8.9 Gb/lane - 71.2 Gb/flowcell
Basecalls: 22 Gb/lane - 176 Gb/flowcell
SRF: 7.53 Gb/lane - 60.42 Gb/flowcell
Filtered Sequences: 1.6 Gb/lane - 12.8 Gb/flowcell
Alignment exports: 1.23 Gb/lane - 14.76 Gb/flowcell
Error reports: 6.47 Gb/lane - 51.76 Gb/flowcell
Minimal Dataset: 8.76 Gb/lane - 70 Gb/flowcell
Everything without images: 66.89 Gb/lane - 535 Gb/flowcell
Everything including images: 166.7 Gb/lane - 1.3 Tb/flowcell