

# Matthew Sooknah

*Computational Biologist & Software Engineer*

1552 East Gate Way, Pleasanton CA 94566  
401-588-2644 — mattsooknah@gmail.com  
mattsooknah.github.io — github.com/mattsooknah

---

## Experience

### 10X Genomics

*Scientist, Computational Biology Group*

**Pleasanton, CA**

*January 2016 – present*

- Developed optimized machine learning methods for profiling and visualization of gene expression across large populations of cells.
- Contributed to processing and analysis of groundbreaking single-cell RNAseq dataset (1.3 million mouse brain cells).
- Helped with design and QC of a custom exome bait set that leverages 10X linked reads to improve haplotype phasing across genes.
- Refined algorithms for short-read alignment and haplotype phasing by using 10X barcodes to reconstruct long DNA molecules.
- Worked with biochemists to perform QC and improve data quality.
- Applied software development best practices to complex data analysis pipelines.

### The Broad Institute of MIT and Harvard

*Software Engineer, Data Sciences & Data Engineering Group*

**Cambridge, MA**

*May 2014 – December 2015*

- Developed pipelines to process petabyte-scale sequencing data produced by the Broad Genomics Platform.
- Contributed to development and support of the Picard and HTSJDK open source toolkits for analyzing sequencing data.
- Wrote high-performance tool to measure incidence of sequencing errors caused by oxidative damage to DNA during preparation of short-read sequencing libraries.
- Helped implement a pipeline for rapid processing and QC of sequencing runs.
- Worked on a backend database and workspace model for FireCloud, a cloud-based platform for analyzing cancer genomics data.
- Developed methods for analyzing gene expression and pathway activity from novel RNAseq-based assays to gain insight into immune cell behavior.

### Nabsys

*Associate Scientist, Algorithms Group*

**Providence, RI**

*June 2013 – May 2014*

- Nabsys develops a microfluidic system that attaches tag molecules to long (~100kb) DNA fragments at sequence-specific sites, then runs them through a nanodetector to produce a map of recognition sites.
- Implemented an improved signal processing pipeline for extracting information about molecules and recognition tags from a noisy electronic readout.
- Prototyped algorithms for assembly, validation and visualization of genomic maps and scaffolds based on Nabsys data.

## Education

**Massachusetts Institute of Technology**  
S.B. Physics, GPA 4.9/5.0

**Cambridge, MA**  
*2009 – 2013*

## Programming Languages

- **Expert:** Python, Java
- **Proficient:** R, Scala, Javascript, bash/awk
- **Familiar:** C/C++, Go, Rust, Scheme, MATLAB

## Analytical Skills

- **Statistics:** hypothesis testing, estimation, error statistics
- **Machine Learning:** clustering, classification, regression, cross-validation, optimization
- **Bioinformatics:** alignment, variant calling, assembly, gene expression, feature annotation

## Publications

- Daniel O’Connell, Raivo Kolde, **Matt Sooknah**, et al. 2016. Simultaneous Pathway Activity Inference and Gene Expression Analysis Using RNA Sequencing. *Cell Systems* 2016; 2(5): 323–334. PMID 27211859.

## Presentations

- “Mapping, processing, and duplicate marking with Picard tools.” BroadE Workshop on GATK Best Practices. Broad Institute, Cambridge, MA. March 2015.

## Poster / Talk Contributions

- Grace Zheng, Jessica Terry, Paul Ryvkin, **Matt Sooknah**, et al. “Single Cell RNA profiling of a Million Neurons by a Massively Parallel and Scalable Droplet Platform”. *Advances in Genome Biology and Technology*. Hollywood Beach, FL. February 2017 (forthcoming).
- Haynes Heaton (presenter), Patrick Marks, **Matt Sooknah**, et al. “Alignment and Variant Calling in Segmental Duplications with Linked-Read Data”. *Genome Informatics*. Wellcome Genome Campus, Hinxton, Cambridge, UK. September 2016.