

# Optimization of Bait Selection for Sequence Capture Arrays

Kevin Silverstein<sup>1</sup>, Kenneth Beckman<sup>2</sup>, Matthew Bower<sup>3</sup>, Kari Bunjer<sup>2</sup>, Matthew Schomaker<sup>4</sup>, Teresa Kemmer<sup>4</sup>, Lisa Schimmenti<sup>5</sup>, Sophia Yohe<sup>6</sup>, Randolph Peterson<sup>6</sup>, Amy Karger<sup>6</sup>, Monika Roychowdhury<sup>6</sup> and Bharat Thyagarajan<sup>6</sup>

<sup>1</sup>Biostatistics and Bioinformatics, Masonic Cancer Center, University of Minnesota, Minneapolis, MN 55455; email: silve023@umn.edu

<sup>2</sup>Biomedical Genomics Center, University of Minnesota, Minneapolis, MN 55455

<sup>3</sup>Division of Genetics and Metabolism, University of Minnesota Medical Center-Fairview, Minneapolis, MN 55455

<sup>4</sup>Molecular Diagnostics Laboratory, University of Minnesota Medical Center-Fairview, Minneapolis, MN 55455

<sup>5</sup>Departments of Pediatrics, Ophthalmology, Genetics, Cell Biology and Development, University of Minnesota, Minneapolis, MN 55455

<sup>6</sup>Laboratory Medicine and Pathology, University of Minnesota, Minneapolis, MN 55455

Next-Generation Sequencing (NGS) technologies will fundamentally alter how clinicians test patients for disease-associated genetic variants. While sequencing costs are falling at a rapid pace, the goal of an affordable full-genome sequence is not likely to be realized in the immediate future. Sequence capture arrays, such as the Agilent SureSelect platform, provide an interim solution allowing NGS technologies to be adapted for specific clinical testing applications. Capture arrays enrich for targeted regions of the genome via hybridization to designed long oligo "baits" for subsequent sequencing by NGS technologies. This offers a cost-effective means for identifying disease-associated variants in genes of interest. To date, oligo bait design has been rather simplistic. Typically, baits are evenly tiled across the region of interest. The final cost of a capture array can be heavily influenced by the total number of bases covered in the reference genome and the number of baits in the design. For large arrays involving many small target exons, the overhang of baits at both ends of these exons can be substantial, nearly doubling the cost of design. We have sought to systematically explore the effects of a variety of layout parameters on the depth of coverage at the boundaries and throughout the targeted exons among ~600 genes. Specifically, we have matched target exons for length, GC content, and problematic bait count (defined as an oligo with GC > 65% or deltaG of folding < 2 SEM below the mean). Within these matched exons (each having N=500), we have systematically compared the default layout against an expanded, compressed, and bait-reduced layout. In another set of matched-exon experiments (each having N=50), we have systematically tested the replacement of problematic baits with nearby improved ones, or enriched their number by 2X or 3X in an attempt to recover difficult-to-capture regions. A comparison of exon mean depth of coverage and fraction of sufficiently covered bases (>10X/20X/30X) across our experimental variables helps us to maximize desired target coverage while minimizing design costs.