



Lea Møller Jagd, Cristina Gamba, and Peter Mouritzen
Samplix ApS, Birkeroed, Denmark

Revealing the structure of a biosynthetic gene cluster in a barley variety using Xdrop®

Summary

- Xdrop enriches specific genomic regions in plant varieties, enabling in-depth sequence elucidation.
- Here, Xdrop is used in a workflow to reveal the structure of a biosynthetic gene cluster in a barley variety with no reference genome.

Introduction

The 5.3 Gb genome of barley (*Hordeum vulgare*) includes several gene clusters involved in the biosynthesis of natural products. Extensive sequencing efforts have generated a barley reference genome and a pan genome,^{1,2} and it is known to contain a large proportion (~80%) of transposable elements with repetitive sequences. Most varieties lack a reference genome.

Xdrop can enrich specific genomic regions for in-depth sequence analysis. Here, we use this method on an important gene cluster in a barley variant.

Experimental setup

The Cer-quc gene cluster of the barley variety PLG-1041 Amsbio, which lacks a reference genome, encodes three genes involved in the production of

β -diketone polyketides. They form part of the wax layer that protects plants against pests and reduces water loss, among its other functions.³

First, we used the Samplix online tool to design three detection sequences within the 100 kb gene cluster based on the reference genome sequence assembly Morex V3. The detection sequences are defined by a single primer pair each. We positioned them as close to the genes as possible (Figure 1). Since the size and sequence of the gene cluster in PLG-1041 Amsbio is unknown, we performed three enrichments. The enriched DNA was sequenced using Oxford Nanopore® technology, yielding 1.9 Gb of data.

Mapping to the Morex V3 reference genome

We mapped the sequence data to the Morex V31 reference genome using Minimap2 (Figure 2). Coverage is high on the gene regions, highlighting the similarities, but few reads map to the intergenic regions, suggesting substantial differences between the two. We performed de novo assembly to characterize the structure in PLG-1041 Amsbio.

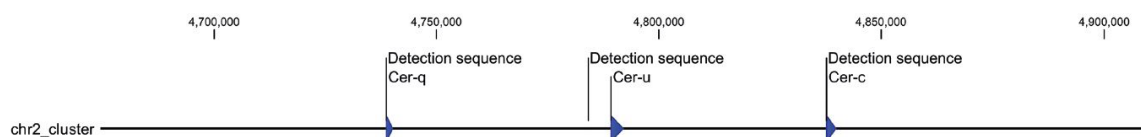


Figure 1. Overview of the Cer-quc gene cluster on chromosome 2H in the More V3 genome assembly. The locations of the three detection sequences are indicated.

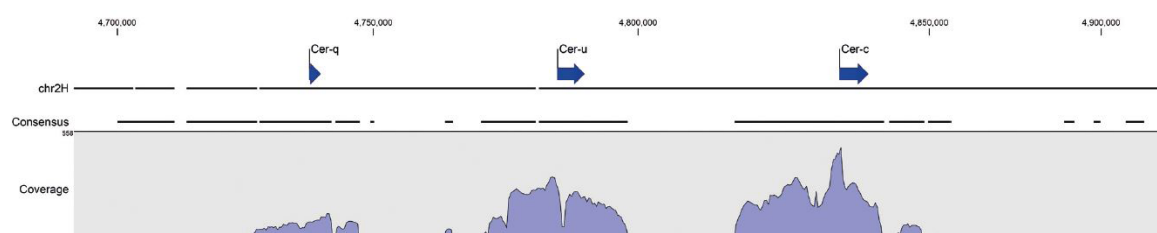


Figure 2. Mapping of the reads obtained in this study to the Morex V3 reference assembly genome using Minimap2. Genes in the cluster show high coverage, whereas limited coverage breadth is observed for the intergenic regions. De novo assembly should be performed to resolve the sequence.

De novo assembly

Prior to de novo assembly, we performed a NECAT⁴ error correction and SACRA⁵ chimera splitting to ensure high quality. Canu⁶ de novo assembly on the resulting reads yielded a 57.5 kb contig (tig00000020) that contains the sequences of Cer-q, Cer-u, and Cer-c. We subsequently mapped the original reads to the contig, finding that the coverage supports the sequence of the assembled region (Figure 3). A LAST⁷ alignment between the contig and the reference showed that the contig contains sequences around the gene regions that are similar to Morex V3, but parts of the Morex V3 reference sequence are not found in the de novo assembly (data available on request).

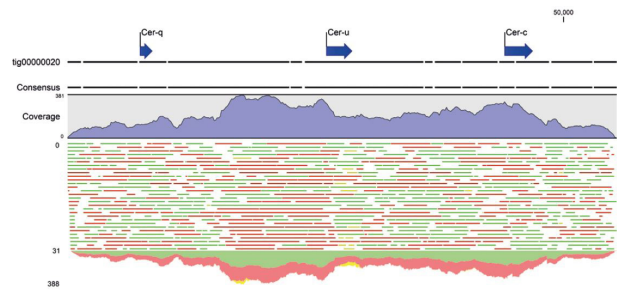


Figure 3. Mapping of the Xdrop enriched reads to contig tig00000020 containing the Cer-q, Cer-u, and Cer-c genes. Red = forward read; green = reverse read; yellow = non-uniquely mapping read.

Analysis of the gene content

Aligning the genes of the reconstructed contig tig00000020 to the genes of Morex V3 demonstrated that the coding sequences are well conserved between the two varieties with some polymorphisms (Figure 4). These polymorphisms are almost exclusively distributed in the introns of the genes, which should not affect the plant phenotype.

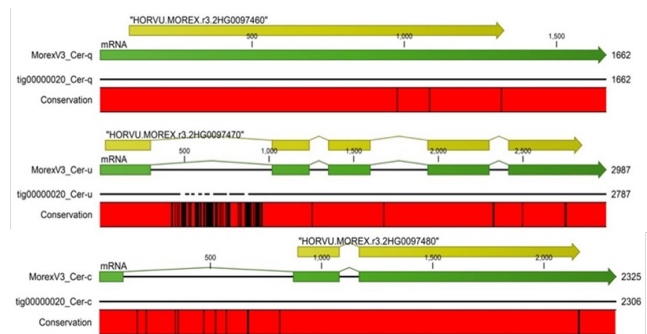


Figure 4. Genes of Morex V3 and contig tig00000020 aligned for comparison. Annotations of cds (yellow) and mRNA (green) were extracted from Morex V3. Conserved nucleotides are shown in red and polymorphisms as black vertical lines.

Conclusion

Xdrop can be used to enrich specific genomic regions in plant varieties, enabling in-depth sequence elucidation, even when only limited knowledge about the variety-specific genomic sequence is available.

For more information about Xdrop products and applications, visit [samplix.com](https://www.samplix.com).

References

- Mascher, M., et al. 2017. A chromosome conformation capture ordered sequence of the barley genome. *Nature* 544: 427. doi: 10.1038/nature22043
- Jayakodi, M., et al. 2020. The barley pan-genome reveals the hidden legacy of mutation breeding. *Nature* 588: 284. doi: 10.1038/s41586-020-2947-8
- Schneider, L.M., et al. 2016. The Cer-cqu gene cluster determines three key players in a β -diketone synthase polyketide pathway synthesizing aliphatics in epicuticular waxes. *Journal of Experimental Botany* 67 (9): 2715–30. doi: 10.1093/jxb/erw105
- Chen, Y. et al. 2021. Efficient assembly of nanopore reads via highly accurate and intact error correction. *Nat Comm.* 12: 60. doi: 10.1038/s41467-020-20236-7
- Kiguchi, Y. et al. 2020. Impact of chimera-less long reads on metagenomics of human gut viromes treated with multiple displacement amplification. Prepr. doi:10.21203/rs.3.rs58640/v1
- Koren, S., et al. 2017. Canu: Scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* 27: 5. doi:10.1101/gr.215087.116
- Kielbasa, S.M., et al. 2011. Adaptive seeds tame genomic sequence comparison. *Genome Res.* 21: 3. doi: 10.1101/gr.113985.110

