Figure1. The minimal sequence difference of each human BCR germline κ light chain gene pair at full sequence length. The colors indicate the number of nucleotide differences between genes (the vertical bar on right gives the scale). The comparison begins at the 3' end of each full-length gene. The numbers in blue circles represent the numbers of mismatched nucleotides in the most similar pairs. These pairs of VH genes cannot be discriminated from each other at almost any length or mutation frequency.
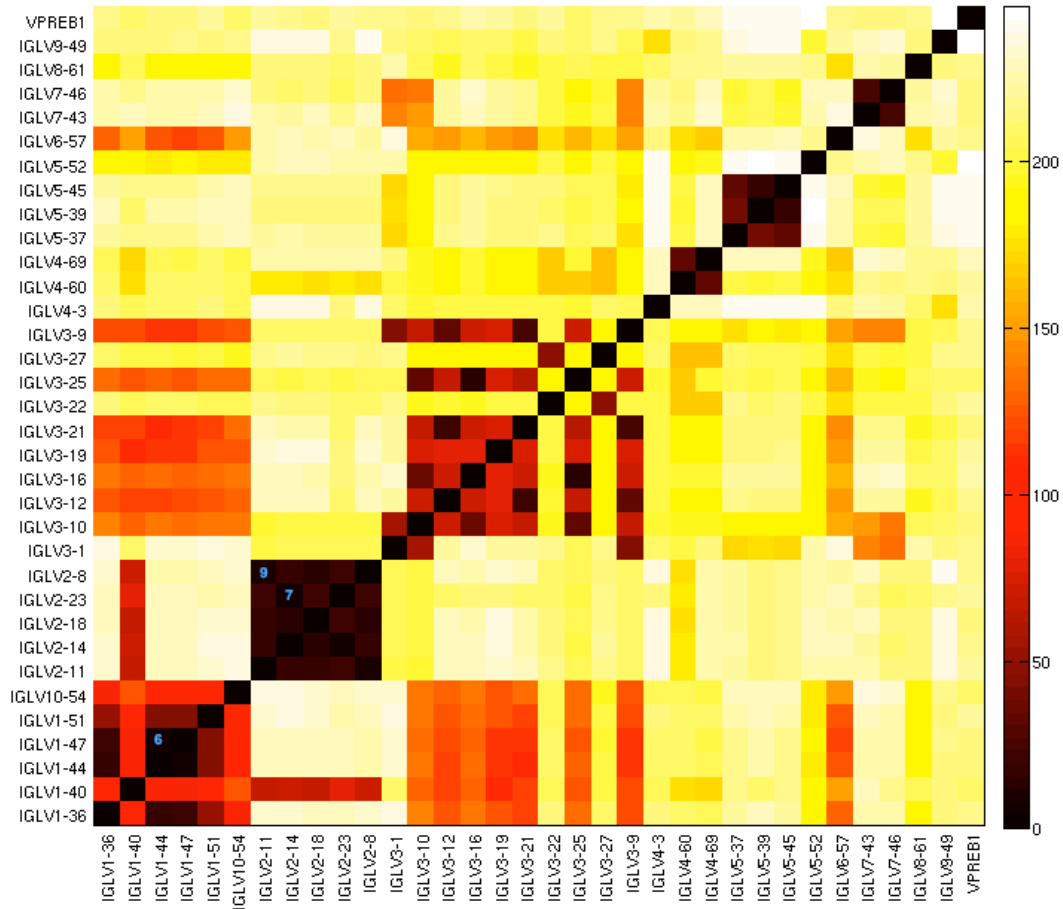
Figure 2. The minimal sequence difference of each human BCR germline λ light chain gene pair at full sequence length. The colors indicate the number of nucleotide differences between genes (the vertical bar on right gives the scale). The comparison begins at the 3' end of each full-length gene. The numbers in blue circles represent the numbers of mismatched nucleotides in the most similar pairs. These pairs of VH genes cannot be discriminated from each other at almost any length or mutation frequency.

Figure 3. The minimal sequence difference of each human TCR germline α chain gene pair at full sequence length. The colors indicate the number of nucleotide differences between genes (the vertical bar on right gives the scale). The comparison begins at the 3' end of each full-length gene. The numbers in blue circles represent the numbers of mismatched nucleotides in the most similar pairs. These pairs of VH genes cannot be discriminated from each other at almost any length or mutation frequency.
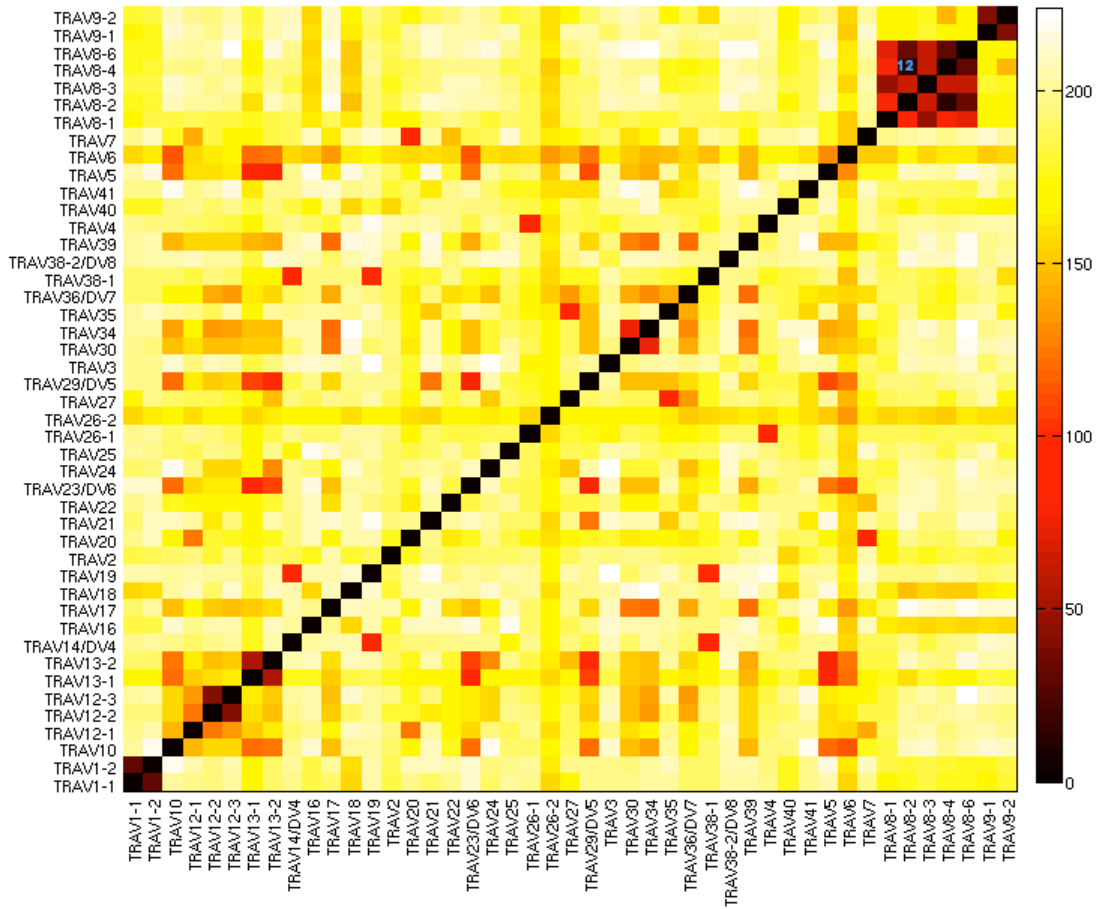
Figure 4. The minimal sequence difference of each human TCR germline β chain gene pair at full sequence length. The colors indicate the number of nucleotide differences between genes (the vertical bar on right gives the scale). The comparison begins at the 3' end of each full-length gene. The numbers in blue circles represent the numbers of mismatched nucleotides in the most similar pairs. These pairs of VH genes cannot be discriminated from each other at almost any length or mutation frequency.
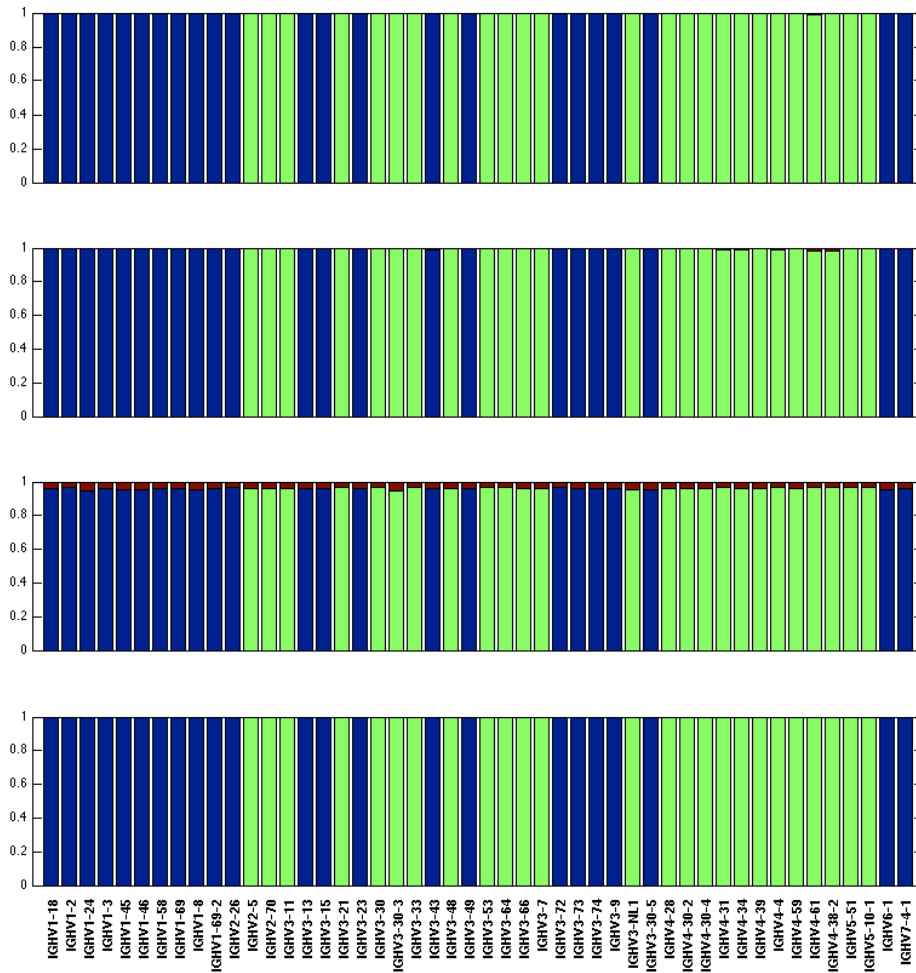
Figure 5. Comparison of human VH identification results using three different sequence identification methods at n=100 nucleotide read length and 0.05 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
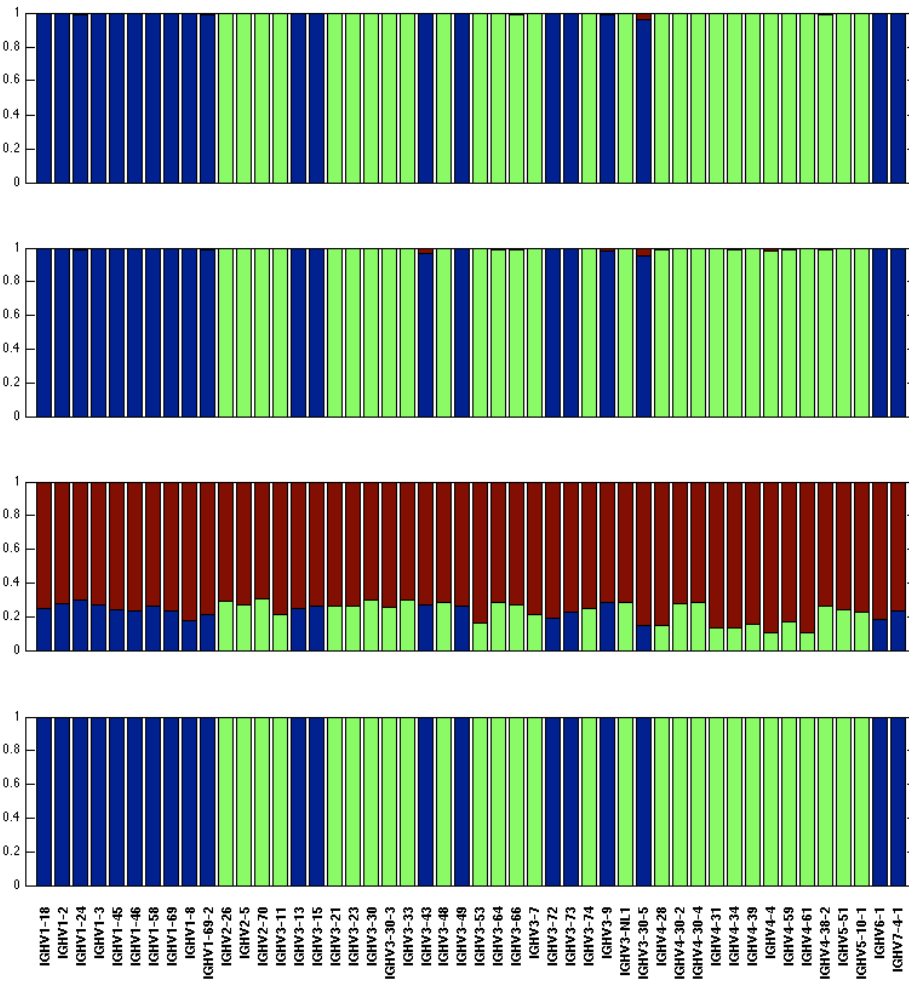
Figure 6. Comparison of human VH identification results using three different sequence identification methods at n=100 nucleotide read length and 0.15 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
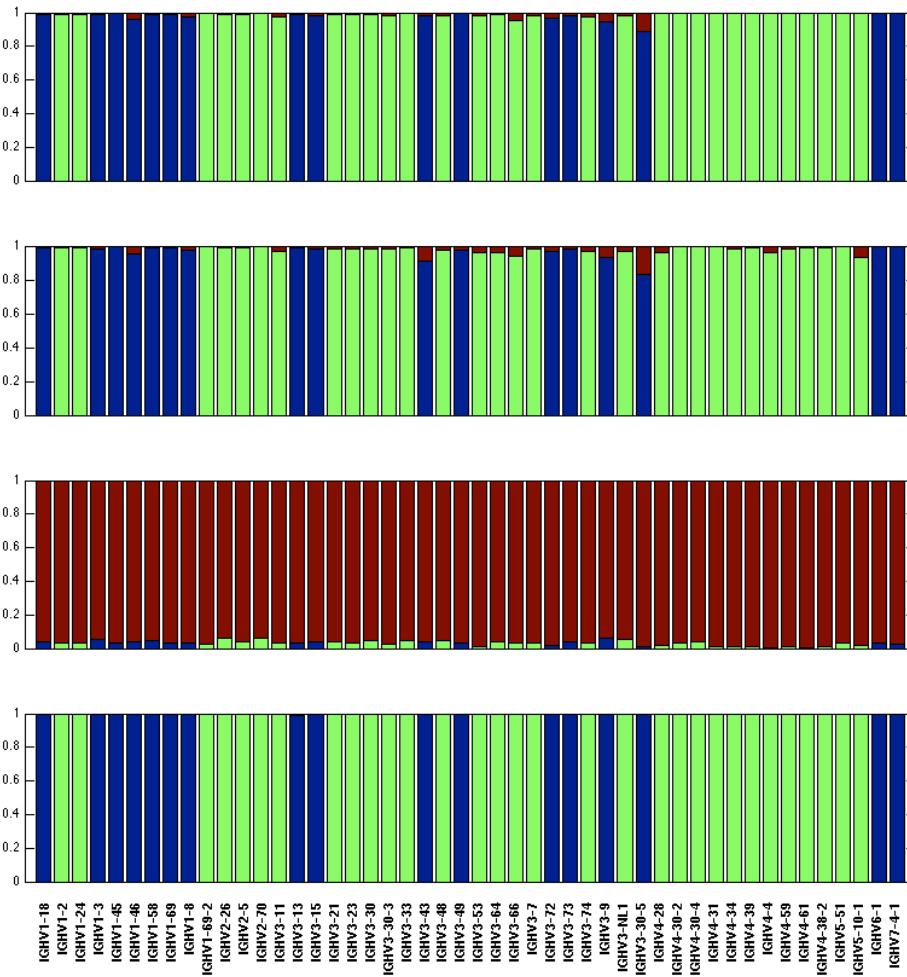
Figure 7. Comparison of human VH identification results using three different sequence identification methods at n=100 nucleotide read length and 0.30 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.

Figure 8. Comparison of human VH identification results using three different sequence identification methods at n=150 nucleotide read length and 0.05 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
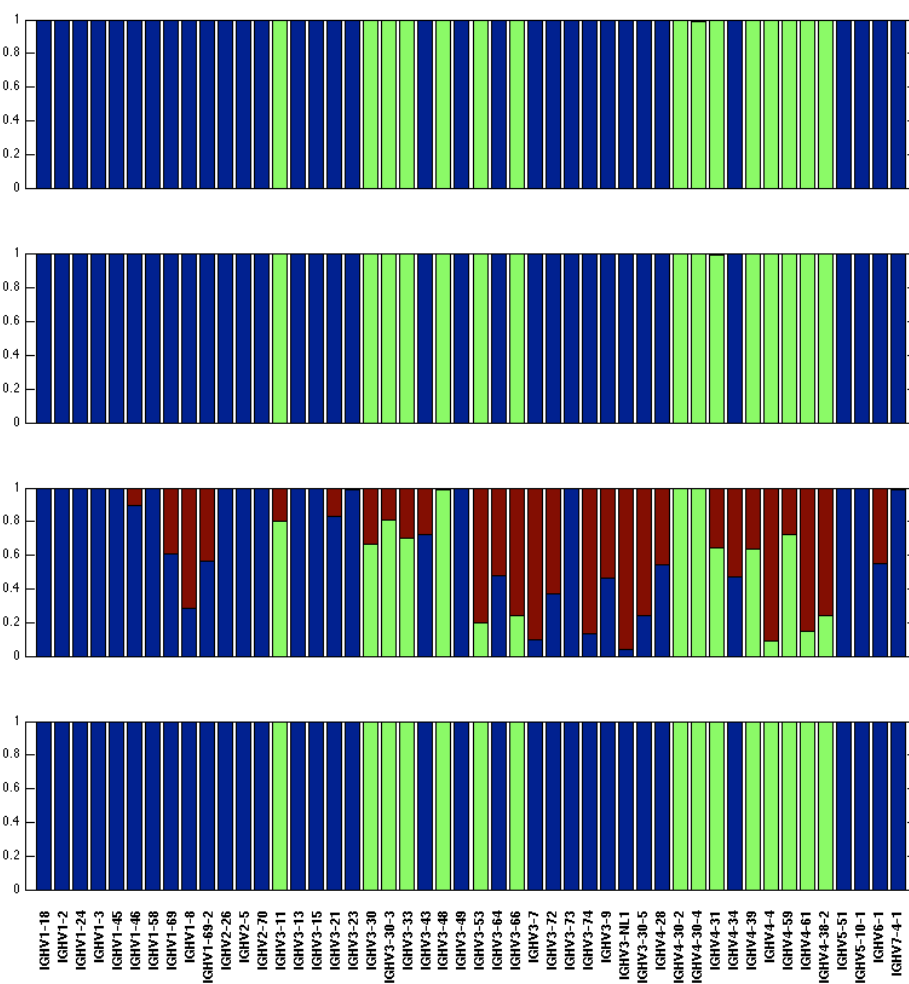
Figure 9. Comparison of human VH identification results using three different sequence identification methods at n=150 nucleotide read length and 0.30 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
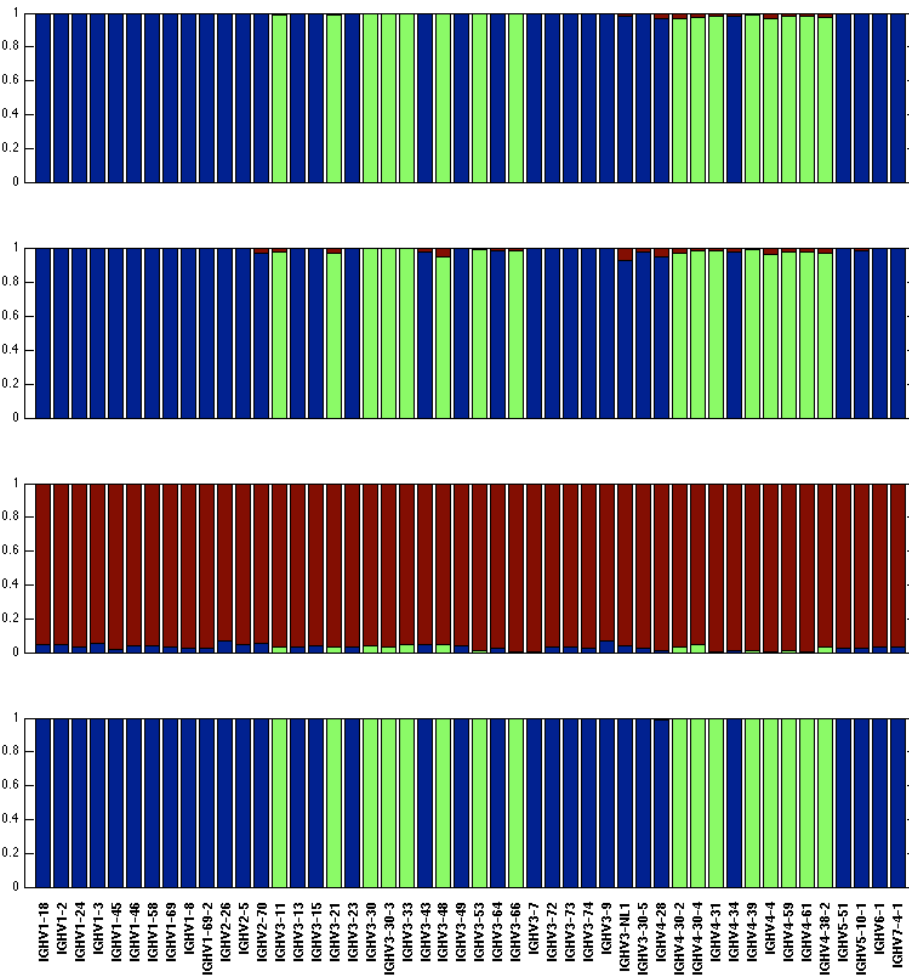
Figure 10. Comparison of human VH identification results using three different sequence identification methods at n=200 nucleotide read length and 0.05 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
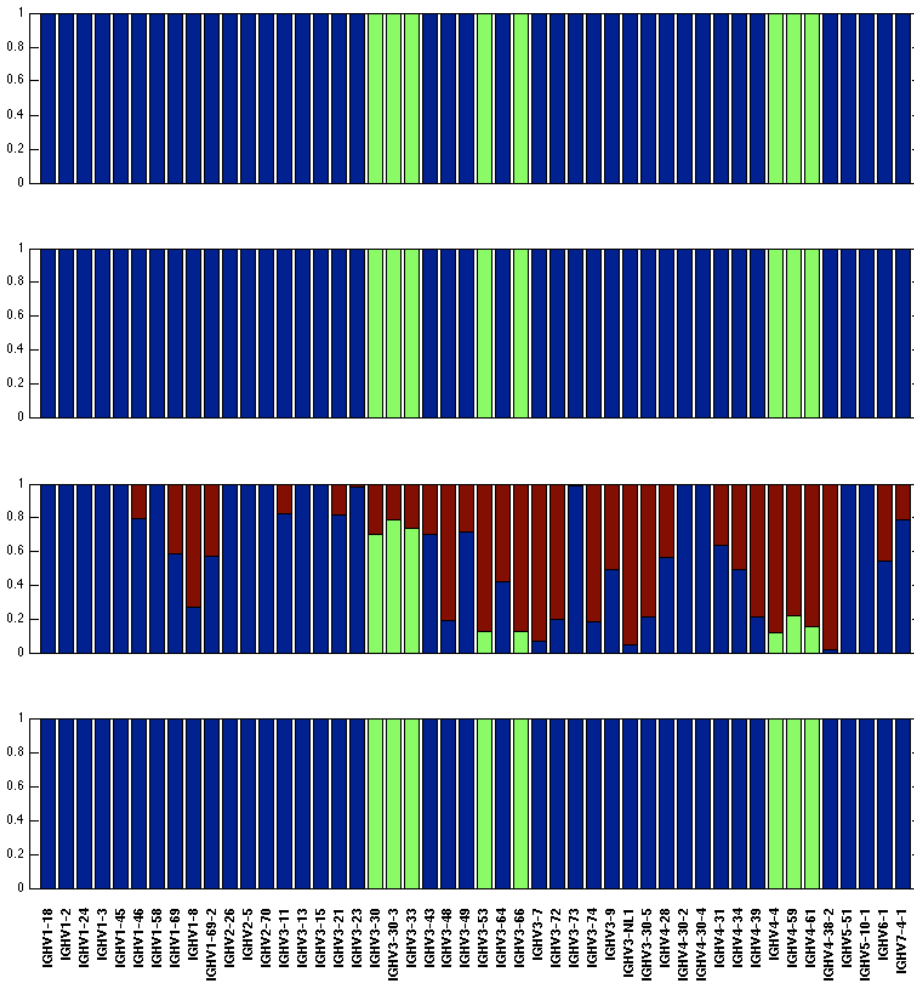
Figure 11. Comparison of human VH identification results using three different sequence identification methods at n=200 nucleotide read length and 0.15 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
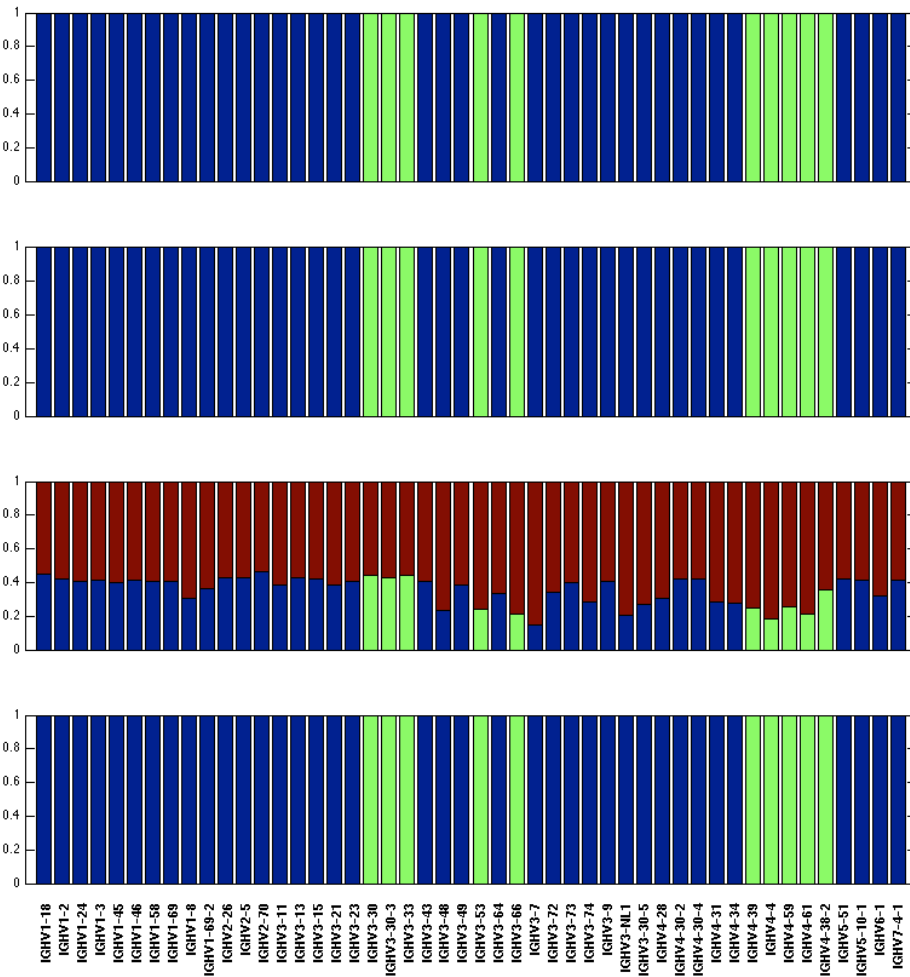
Figure 12. Comparison of human VH identification results using three different sequence identification methods at n=200 nucleotide read length and 0.30 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
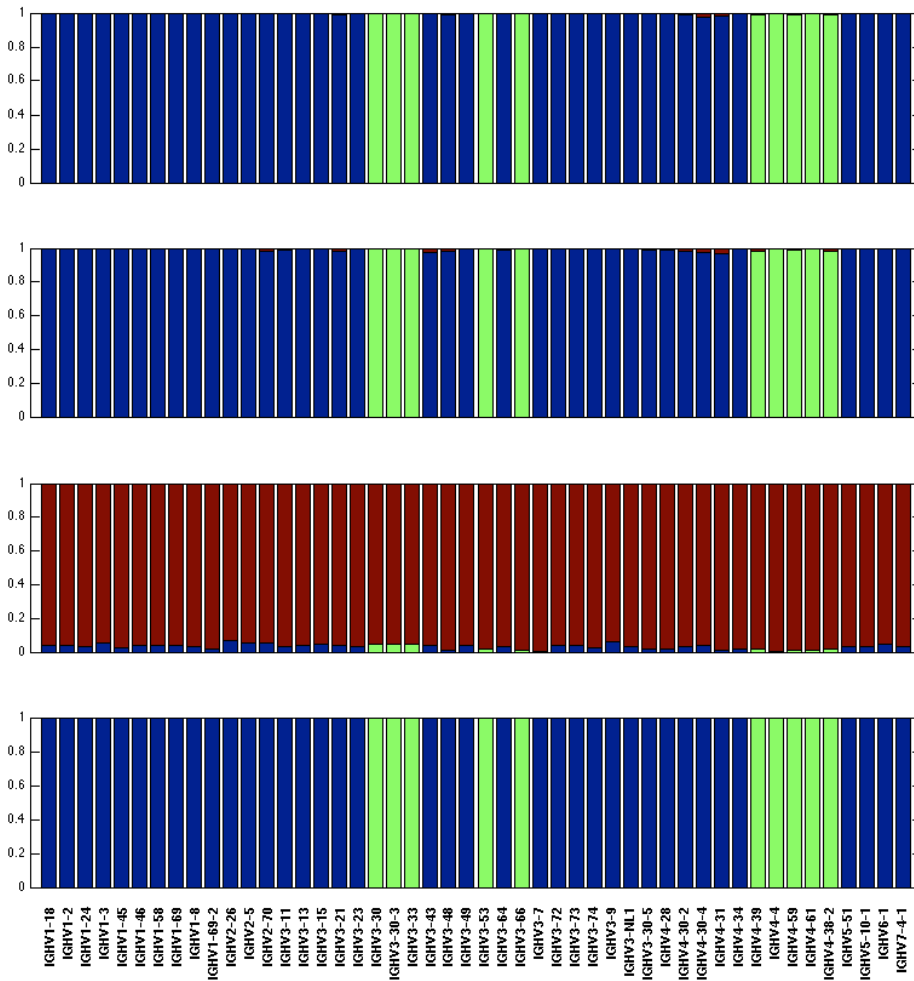
Figure 13. Comparison of human VH identification results using three different sequence identification methods at full read length and 0.05 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
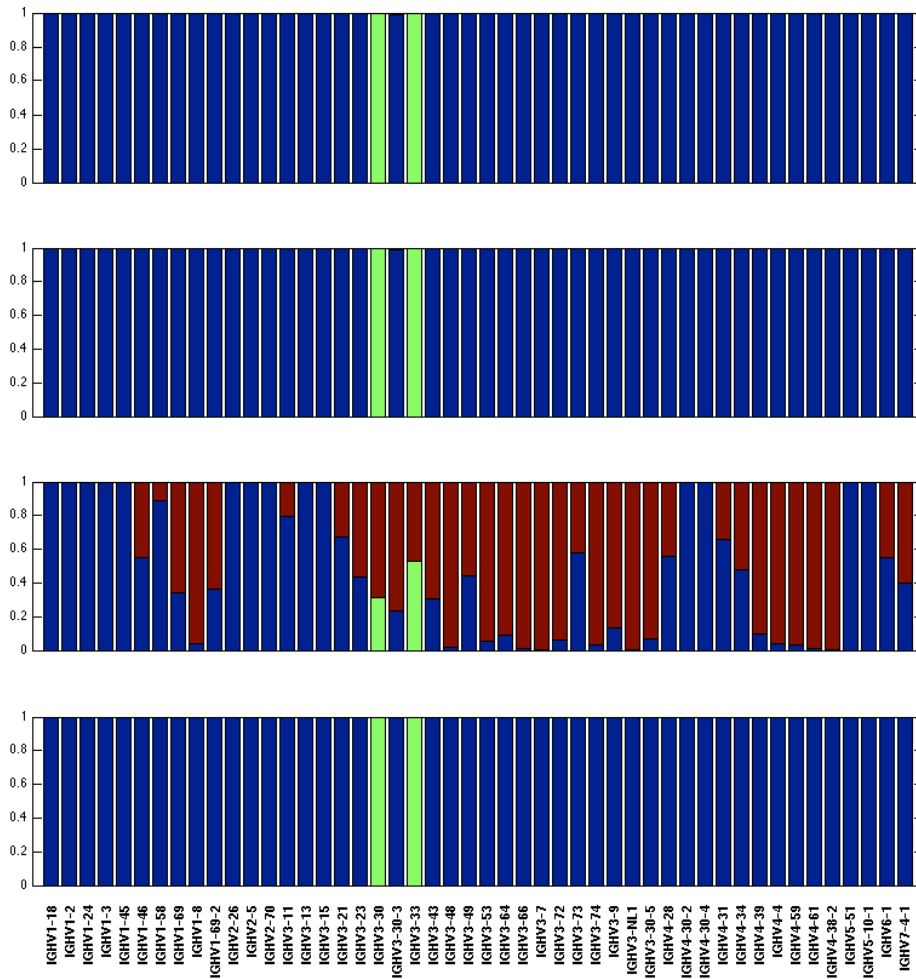
Figure 14. Comparison of human VH identification results using three different sequence identification methods at full read length and 0.15 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.
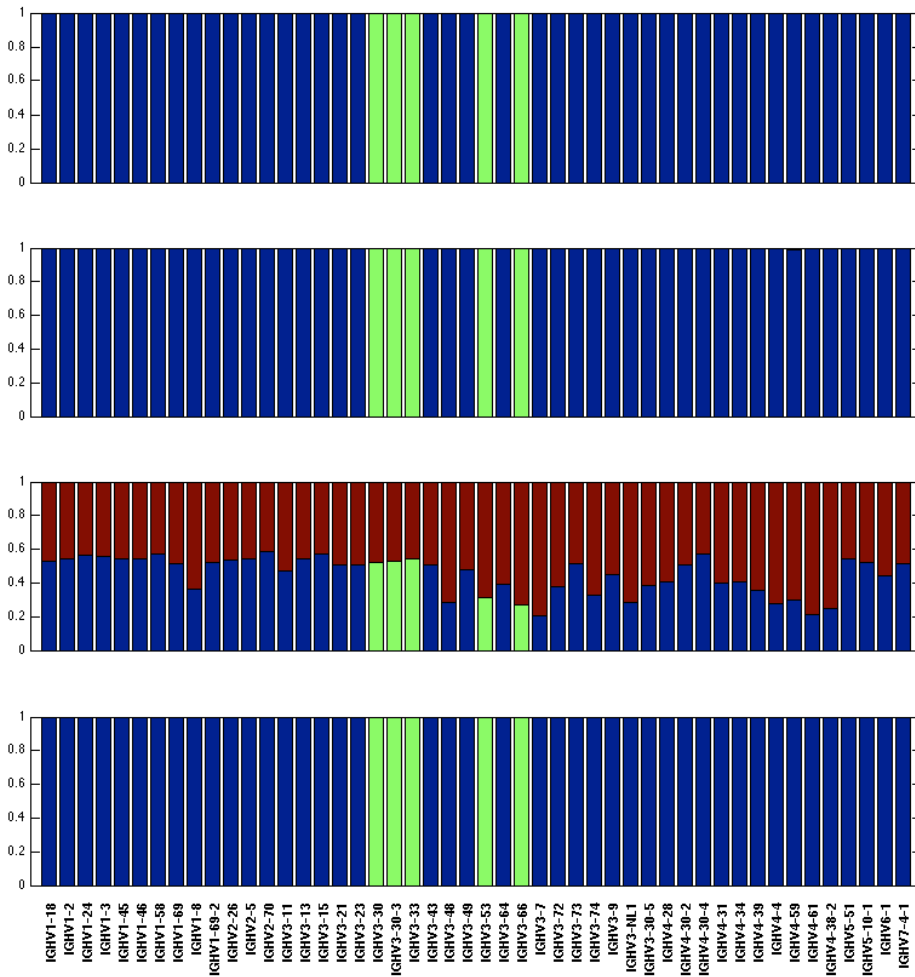
Figure 15. Comparison of human VH identification results using three different sequence identification methods at full read length and 0.30 mutation frequency. From top to bottom (A) The Anchor method (this paper); (B) High V-Quest; (C) IgBLAST using the default word length (9 characters); (D) IgBLAST using the minimal word length (4 characters). In all cases the blue color represents the fraction of correctly identified sequences in which a VH gene is uniquely identifiable; green is the fraction of identifications when a gene is not distinguishable from at least one other gene (using the calculation described in the Calcualtion and Methods section) and either the real gene or the one we predict to confuse it are identified; red is the fraction of incorrect identifications.