

Subset seed extension to protein BLAST

Anna Gambin Sławomir Lasota Laurent Noé Michał Startek,
Maciej Sykulski and Gregory Kucherov

September 21, 2009

Appendix

Evolution of seed alphabets and seeds

The fitness functions considered in the paper are as follows:

$$\text{fitness}_1(M) = \frac{\text{sens}^F(M)}{\text{sens}^B(M)}. \quad (1)$$

$$\text{fitness}_2(M) = \begin{cases} 0 & \text{if } \text{sens}^B(M) > c \\ \text{sens}^F(M) & \text{otherwise} \end{cases}$$

$$\text{fitness}_3(M) = \begin{cases} \text{sens}^F(M) & \text{if } \text{sens}^F(M) < c \\ \frac{\text{sens}^F(M)}{\text{sens}^B(M)} & \text{otherwise} \end{cases}$$

$$\text{fitness}_4(M) = \frac{\text{sens}^F(M)}{\text{sens}^B(M) + c}$$

It turned out that $\text{fitness}_1(M)$ yields unsatisfactory results – the evolution just results in a smallest multiple seed possible, with minuscule foreground and background sensitivity (see Figure 1).

Fitness functions $\text{fitness}_2(M)$, $\text{fitness}_3(M)$ and $\text{fitness}_4(M)$ reflect the trade-off between foreground sensitivity and background sensitivity. The resulting path of evolution for $\text{fitness}_2(M)$ is shown on Figure 2, in this case for $c = 0,002195$. Figure 3 presents the path of evolution for $\text{fitness}_3(M)$ with $c = 0.2$ and Figure 4 for $\text{fitness}_4(M)$ and $c = 0.001$.

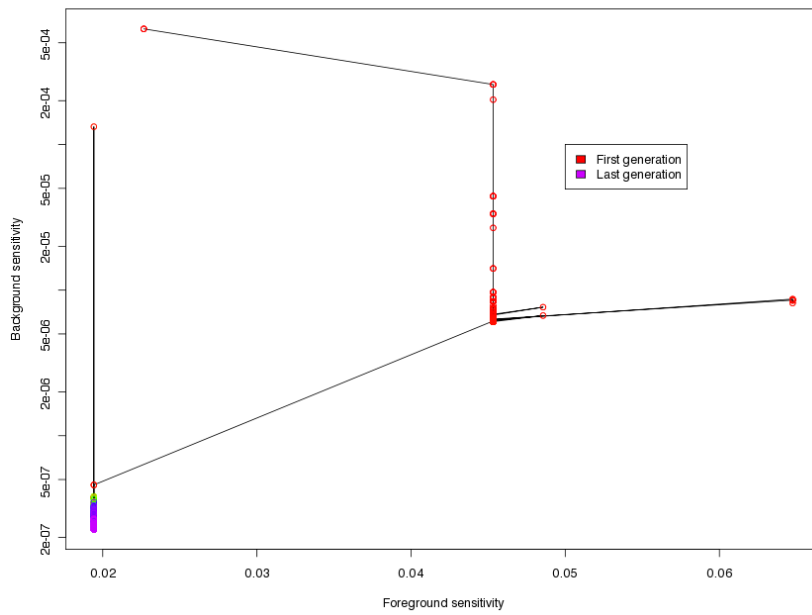


Figure 1: Path of evolution for fitness function $fitness_1$.

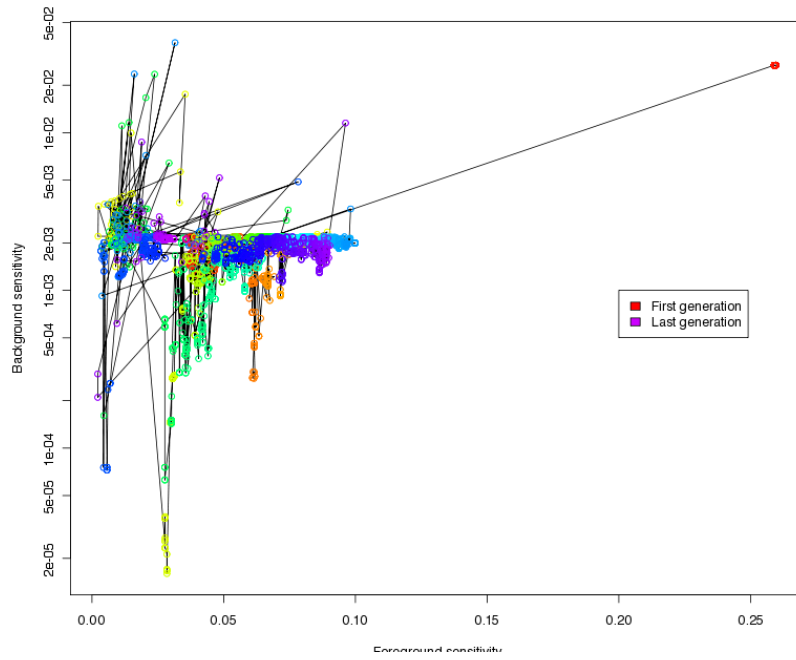


Figure 2: Path of evolution for fitness function $fitness_2$.

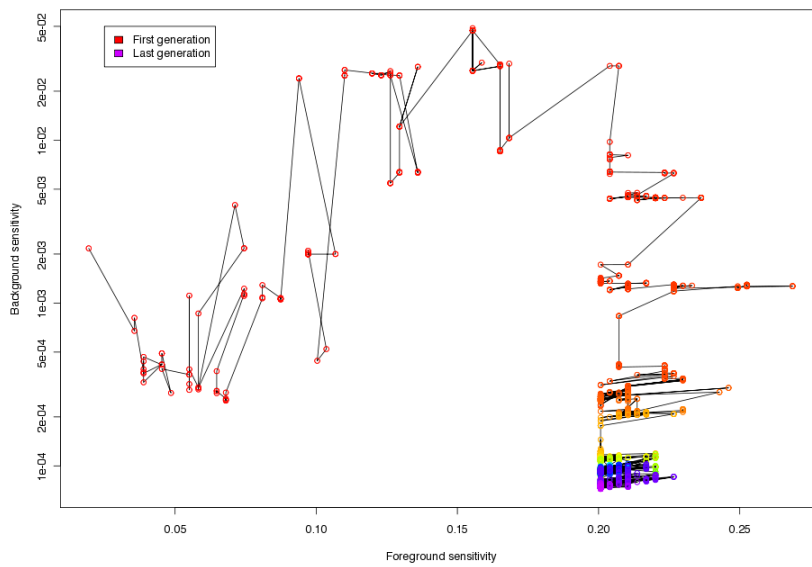


Figure 3: Path of evolution for fitness function $fitness_3$.

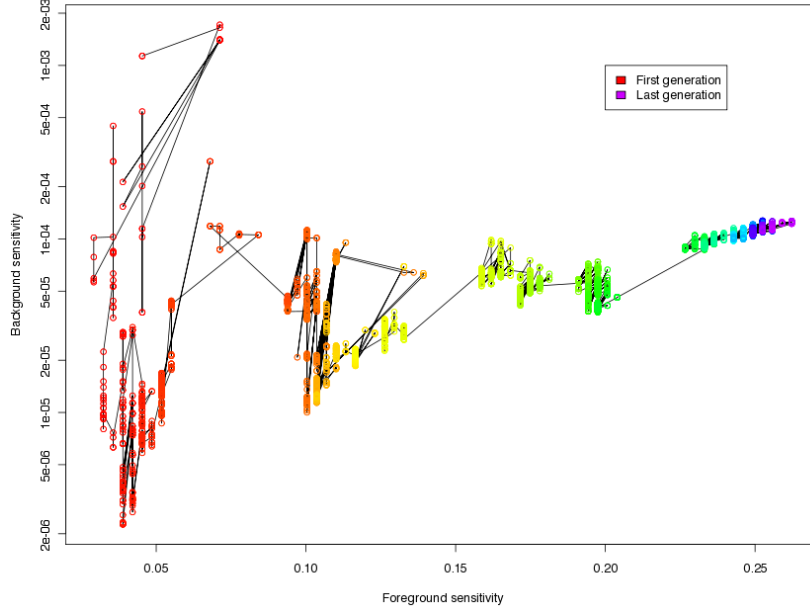


Figure 4: Path of evolution for fitness function $fitness_4$.

Alphabet of the multiple seed used for experiments

The alphabet contains 8 letters. We provide its representation as a hierarchy of nested equivalence relations over amino acid alphabets. The hierarchical grouping of amino acids is also presented.

- : {CFYWMLIVGPATSNHQEDRK}
- 1 : {RKQED, IVLFM, A, S, Y, T, G, N, H, C, P, W}
- 2 : {RKQED, IVL, F, M, A, S, Y, T, G, N, H, C, P, W}
- 3 : {RKQ, ED, IVL, F, M, A, S, Y, T, G, N, H, C, P, W}
- 4 : {RK, Q, ED, IVL, F, M, A, S, Y, T, G, N, H, C, P, W}
- 5 : {RK, Q, E, D, IV, L, F, M, A, S, Y, T, G, N, H, C, P, W}
- 6 : {R, K, Q, E, D, IV, L, F, M, A, S, Y, T, G, N, H, C, P, W}
- # : {R, K, Q, E, D, I, V, L, F, M, A, S, Y, T, G, N, H, C, P, W}

(2)

