

File S1

Supporting Information

Power and False Positive Rate. In order to evaluate the power of SFselect and XP-SFselect to detect positive selection as compared to other neutrality tests, we applied these tests to several datasets simulated under different model parameters. For a given test on a given dataset, the power at 5% false positive rate (FPR) was estimated as the fraction of test-statistic values exceeding a set threshold when applied to the selected samples. The threshold was set to the top 5% of the null distribution, obtained by applying the test to neutral samples. For cross-populations tests (including XP-SFselect) we used the same procedure, only applying the test to selected vs. neutral samples, while the null was obtained by applying the test to neutral1 vs. neutral2 samples.

SVM implementation details. We used a linear (dot product) kernel function SVM. Linear kernels have two important advantages. First, because feature-weights learned by a linear SVM represent a maximum-margin separating hyperplane of the training data in the problem space (rather than in a higher dimensional space), they correspond to the relative importance of features in separating the training data, making the trained SVM easily interpretable. Secondly, normalization of the training and testing data is done in the input space, without the need for complicated normalization of the kernel function itself (Graf *et al.* 2003).

The SVM implementation we used was from the LIBSVM library (Chang and Lin 2011), packaged in the python library scikit-learn (Pedregosa *et al.* 2011). For the parameter-specific SVMs, where we lacked sufficient simulated data to hold the test data out of training, we report power as the mean over 50-fold cross validation. For the general two-stage SVM (SFselect and XP-SFselect), testing and training were done on completely separate datasets.

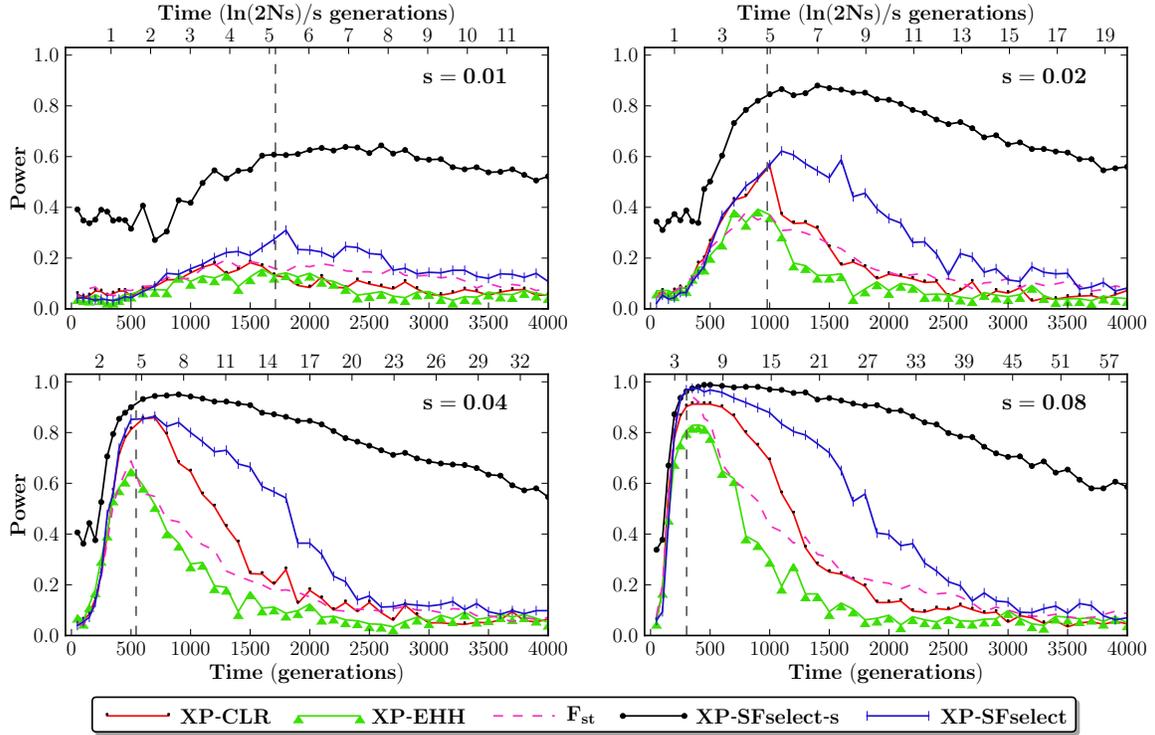


Figure S1: Power (0.05 FPR) of the cross-population SVM test compared to other cross-population tests of neutrality. Shown across selection pressures $s \in [0.01, 0.08]$ and times $\tau \in [0, 4000]$. The (black) line labelled ‘XP-SFselect-s’ shows power when assuming knowledge of the selection coefficient and the time (τ and s , respectively). The (blue) line labeled ‘XP-SFselect’ shows power when no prior knowledge of (s, τ) is assumed. Time is shown in generations (bottom axes), and $\ln(2Ns)/s$ generations (top axes). The dashed vertical lines (grey) show the mean time to fixation of the beneficial allele, which occurs at $\approx 5 \ln(2Ns)/s$ generations.

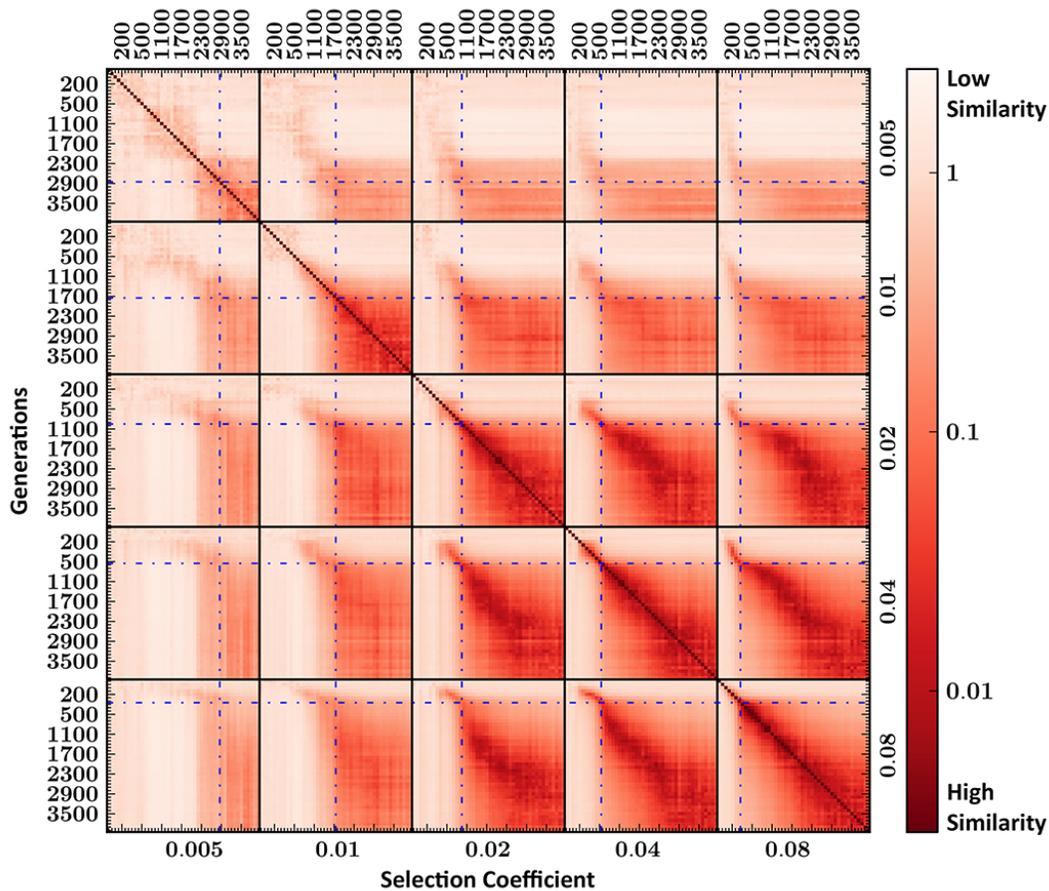


Figure S2: Pairwise cosine distance between 200 SVMs trained on cross-population data (matrices of the XP-SFS scaled to 8×8 frequency bins, and vectorized). The data was simulated under different selection pressures $s \in [0.005, 0.08]$, and sampled at different times under selection $\tau \in [0, 4000]$ generations. Selection pressure boundaries are denoted by black lines, and mean time to fixation for each pressure is denoted by dashed blue lines. We observe two main similarity blocks at each selection pressure, corresponding to "near fixation" and "post-fixation" of the beneficial allele. The stronger the selection pressure (e.g., bottom right) the earlier and shorter the near-fixation stage, and vice versa.

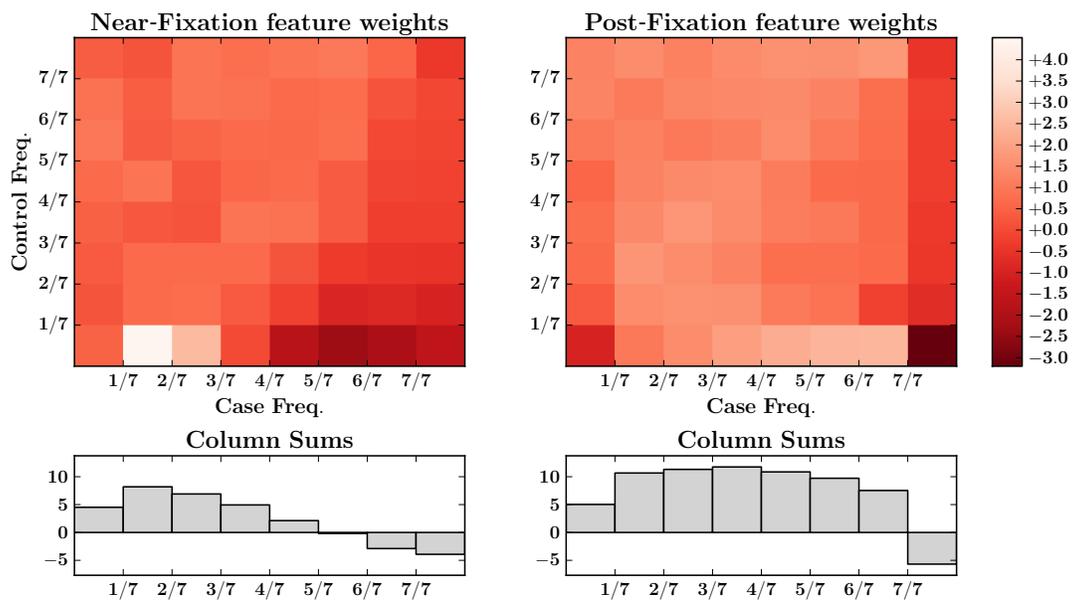


Figure S3: Feature weights learned from the XP-SFS on data corresponding to the two observed regimes of selection: (A) near-fixation, and (B) post-fixation. Minor allele frequencies were distributed to 8×8 bins, where the rightmost column (top row) was dedicated to alleles fixed in the selected (neutral) population. Decision function constants were $\beta_0 = -0.80$, and $\beta_0 = -0.56$ for the near-fixation, and post-fixation SVMs, respectively.

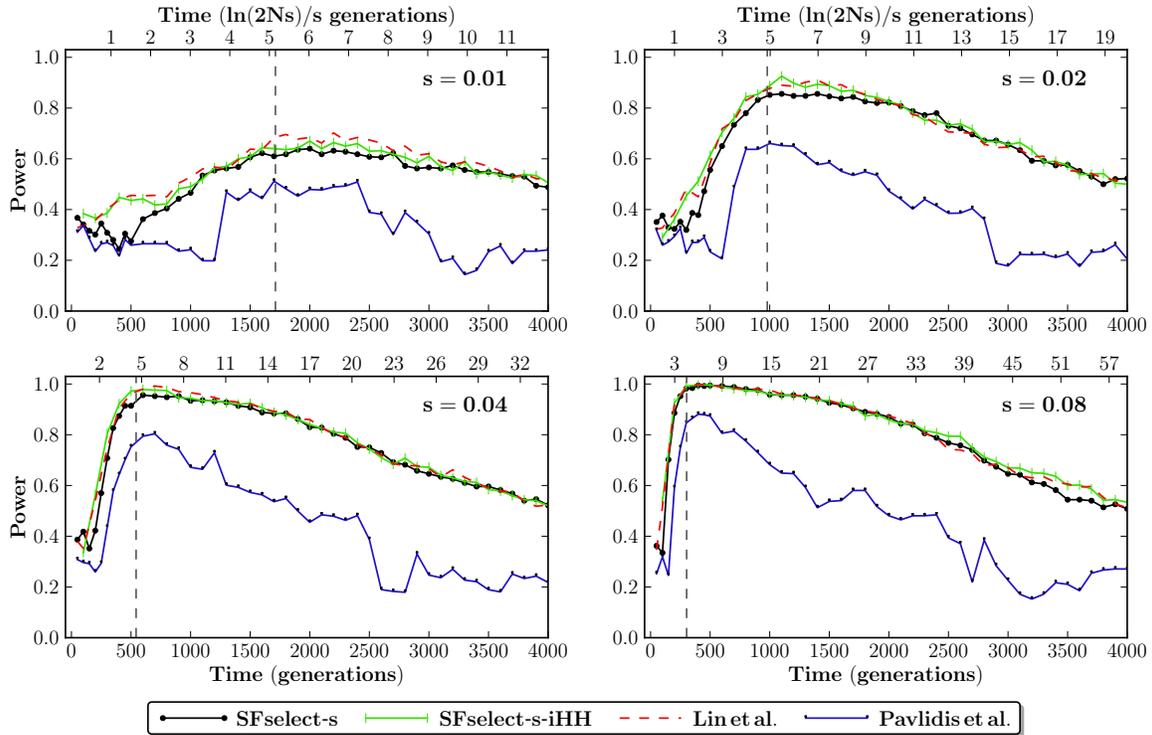


Figure S4: Power (0.05 FPR) of neutrality tests based on supervised learning. The line labelled ‘SFselect- s ’ shows power of the regular parameter-specific SVMs, while the line labelled ‘SFselect- s -iHH’ shows power when including the iHH features described in Lin *et al.* (2011). Shown for selection pressures $s \in [0.01, 0.08]$ and times $\tau \in [0, 4000]$, with time in generations (bottom axes), and $\ln(2Ns)/s$ generations (top axes). The dashed vertical lines (grey) show the mean time to fixation of the beneficial allele, which occurs at $\approx 5 \ln(2Ns)/s$ generations.

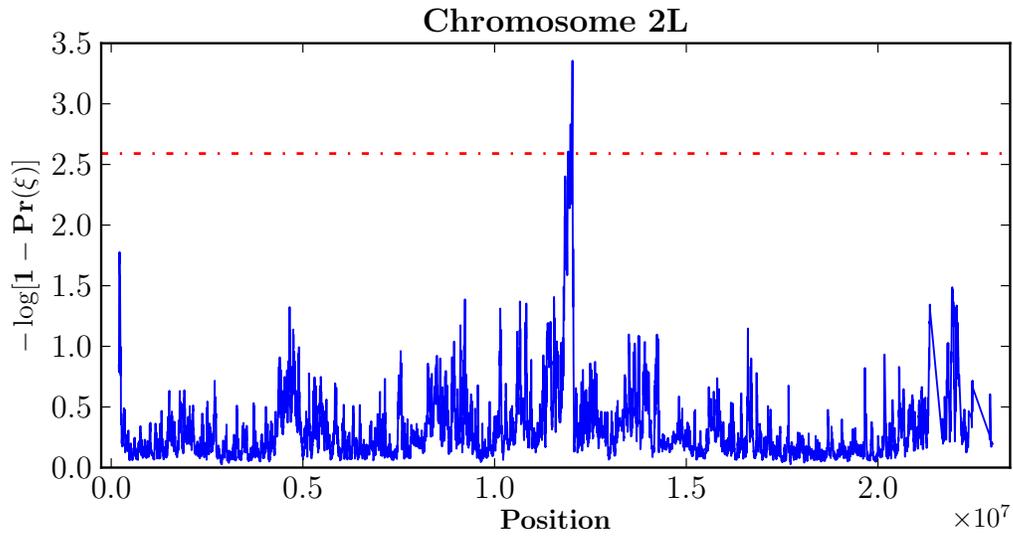


Figure S5: XP-SFselect values on Drosophila chromosome 2L.

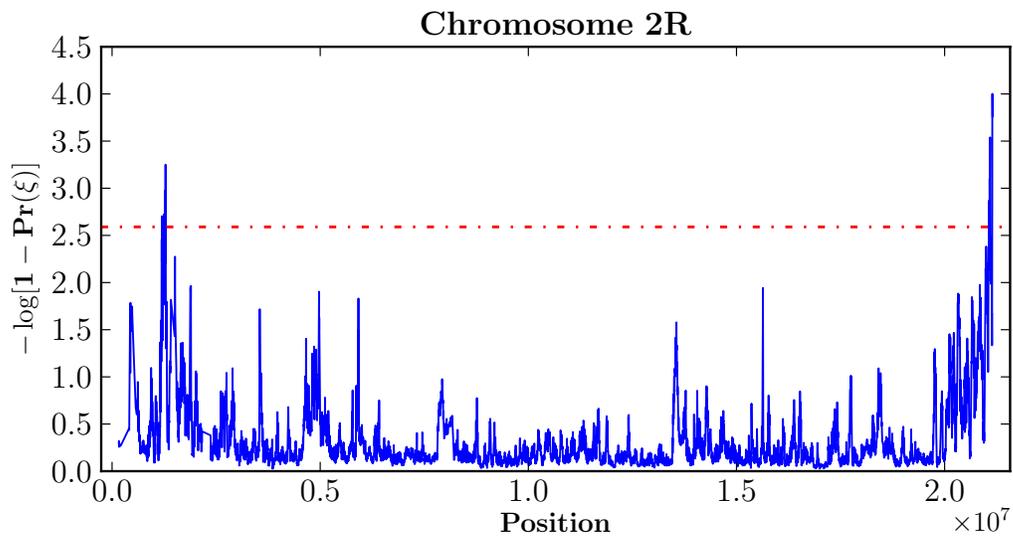


Figure S6: XP-SFselect values on Drosophila chromosome 2R.

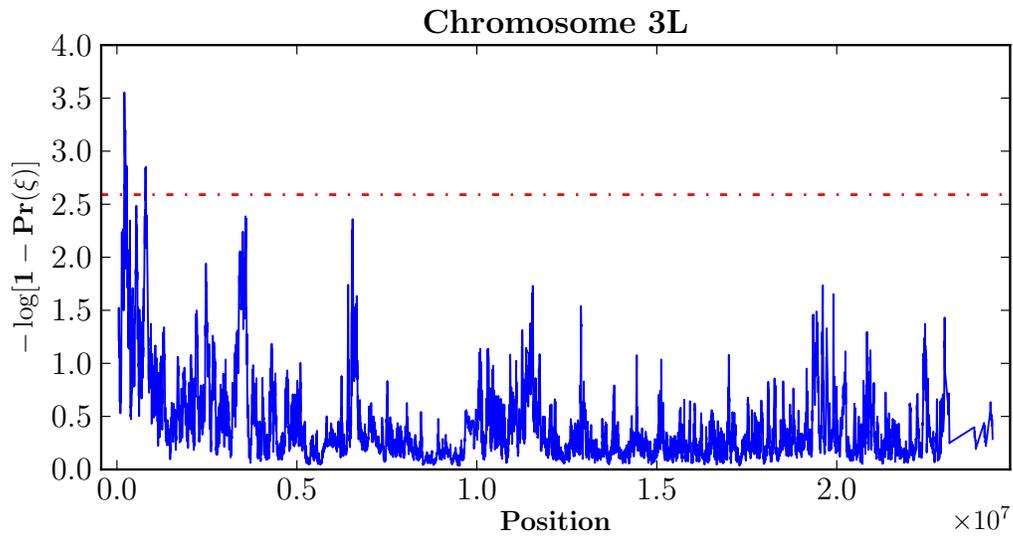


Figure S7: XP-SFselect values on *Drosophila* chromosome 3L.

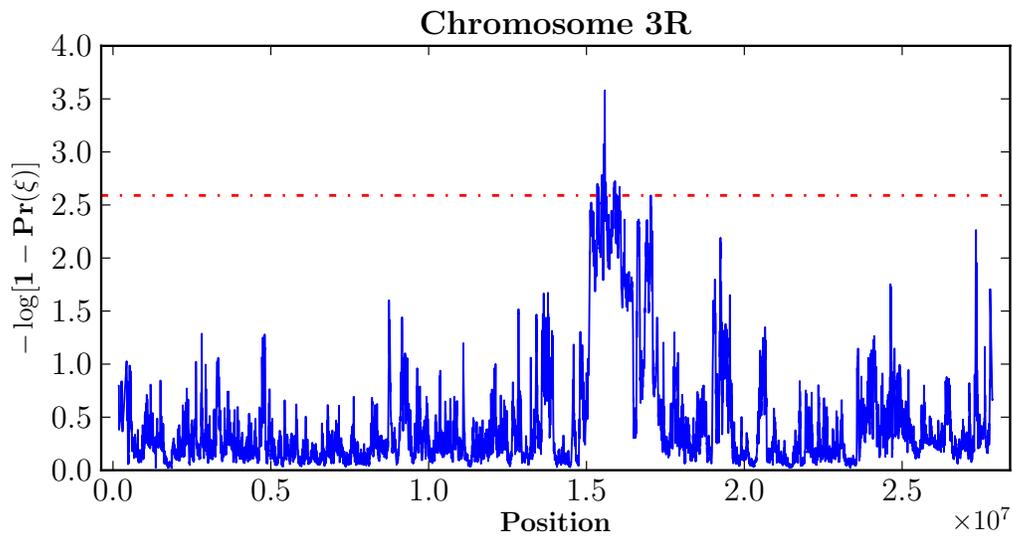


Figure S8: XP-SFselect values on *Drosophila* chromosome 3R.

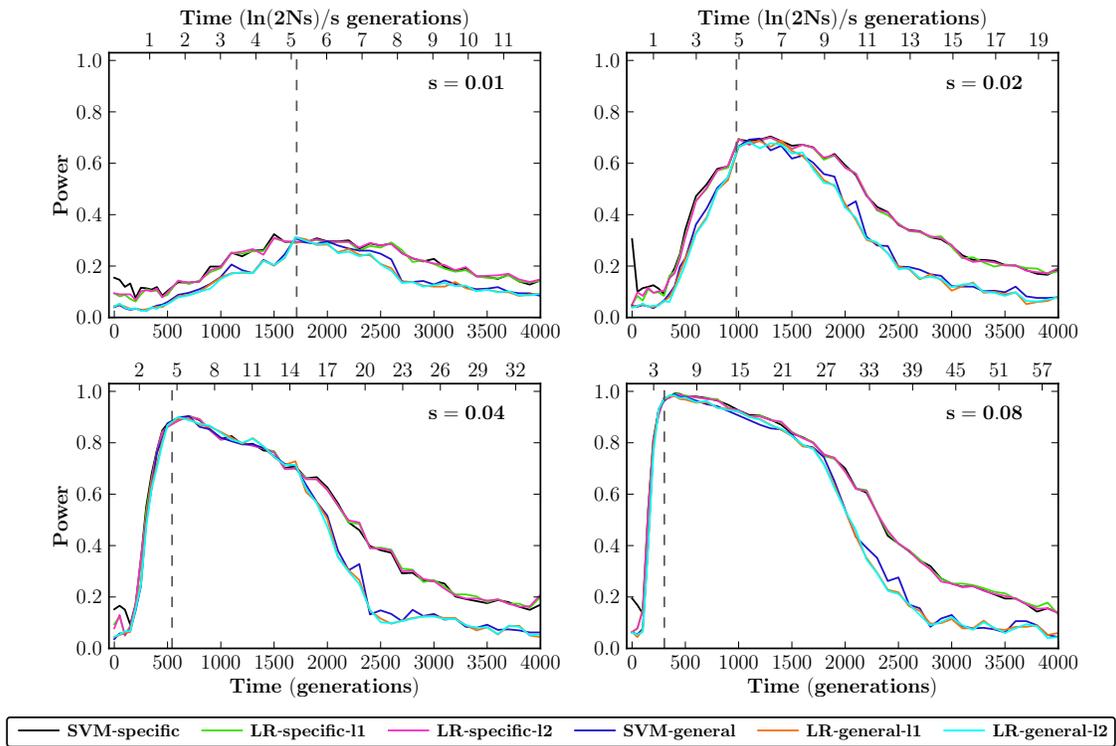


Figure S9: Power of SFselect using SVM and logistic regression, at different times and selection pressures. Performance appears nearly identical regardless of the underlying classification method. In the legend, '11' and '12' refer to the regularization term used with logistic regression.

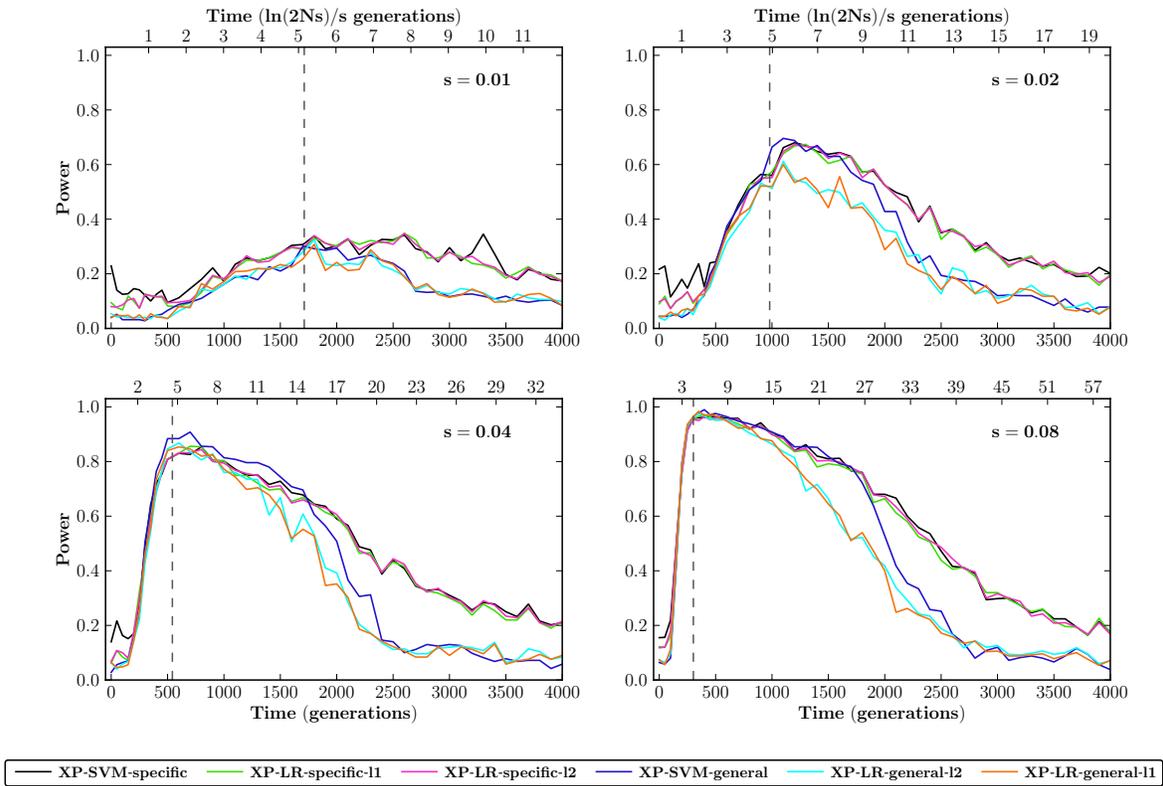


Figure S10: Power of XP-SFselect using SVM and logistic regression, at different times and selection pressures. We observe a marked decrease in power (cyan and orange (LR) compared to blue (SVM)) with logistic regression. In the legend, 'l1' and 'l2' refer to the regularization term used with logistic regression.

Table S1: List of significant regions under XP-SFselect for the fly hypoxia experiments described in Zhou *et al.* (2011).

Chr	Region	XP-SFselect
2L	11895542-12055542	3.36
2R	170962-1308962	3.25
2R	1040962-21174962	4.00
3L	175762-301762	3.55
3L	763762-833762	2.85
3R	15318233-15642233*	3.58
3R	15846233-16076233*	2.73
3R	17014233-17064233	2.59
X	378615-440615	2.78
X	676615-728615*	2.60
X	1420615-1480615	2.70
X	2046615-2122615	4.77
X	2630615-2758615	3.91
X	2872615-3444615	3.89
X	4818615-4892615	3.60
X	12996615-13374615*	4.02
X	15092615-15160615	2.84
X	16110615-16160615*	2.62
X	16276615-16488615*	4.86
X	18154615-18248615	2.87
X	18564615-18686615*	2.60
X	18838615-18930615*	3.08
X	19092615-19358615*	3.74
X	20504615-20986615*	4.18
X	22064615-22412615*	4.24

*Shared with S_f

Table S2: The top 40 non-overlapping regions identified genome-wide by XP-SFselect.

Chr	Position (Mb)	Max XP-SFselect	Genes	Study
X	66.10-66.56	4.38657		
12	88.24-88.36	4.38258		
X	99.00-99.16	4.34082	LOC442459	Frazer <i>et al.</i> (2007)
8	52.67-52.82	4.31338	PXDNL, PCMTD1	(Frazer <i>et al.</i> 2007)
X	35.27-35.38	4.27039		
12	123.61-123.78	4.20905	MPHOSPH9, C12orf65, CDK2AP1, SBNO1	
12	88.90-89.00	4.20736	KITLG	(Pickrell <i>et al.</i> 2009)
4	148.54-148.79	4.19501	TMEM184C, PRMT10, ARHGAP10	
10	100.78-100.94	4.19111	HPSE2	
10	31.47-31.55	4.14863		(Chen <i>et al.</i> 2010)
X	110.08-110.37	4.13684	PAK3	(Sabeti <i>et al.</i> 2007); (Frazer <i>et al.</i> 2007)
2	13.69-13.90	4.12967		
11	105.99-106.22	4.11825		
X	80.24-80.38	4.10921	HMGN5	
13	71.98-72.12	4.10127	DACH1	
4	52.88-53.14	4.09292	LRRC66, SGCB, SPATA18	
2	150.39-150.49	4.07069	MMADHC	
15	44.29-44.39	4.05108	FRMD5	
1	142.66-142.87	4.04074		
11	40.22-40.32	4.02215	LRRC4C	
16	15.14-15.30	4.01629	NTAN1, RRN3, MIR3180-4	
2	97.68-97.85	3.99585	FAHD2B, ANKRD36	
4	159.35-159.44	3.9884		
2	104.76-104.83	3.97891		
17	73.30-73.44	3.96581	GRB2, MIR3678	
20	60.66-60.73	3.93865	LSM14B, PSMA7, SS18L1	
4	41.96-42.11	3.93681	TMEM33, DCAF4L1, SLC30A9	
15	28.19-28.27	3.91923	OCA2	(Chen <i>et al.</i> 2010)
1	158.15-158.24	3.89804	CD1D, CD1A	
13	41.39-41.54	3.8952	SUGT1P3, ELF1	
1	100.67-100.77	3.88985	DBT, RTCD1, MIR553	
X	65.54-65.91	3.87444	EDA2R	
17	53.79-53.87	3.87161	TMEM100, PCTP	
18	30.40-30.58	3.86989	C18orf34	
1	248.07-248.16	3.86911	OR2T8, OR2L13, OR2L81, OR2AK2, OR2L1P	
16	79.80-79.88	3.86909		(Chen <i>et al.</i> 2010); (Frazer <i>et al.</i> 2007)
X	108.00-108.15	3.82083		
18	15.04-15.15	3.81846		(Frazer <i>et al.</i> 2007)
2	167.50-167.60	3.81693		
X	74.42-74.72	3.80593	UPRT, ZDHHC15	

The right-most column specifies the studies, if any, in which the corresponding regions were reported as showing signal of selection.

Table S3: Potentially damaging SNPs found in regions with strong evidence of non-neutral evolution.

Chr	Position	rsID	AA	SIFT	Gene	ENSEMBL	CEU	YRI
1	11090916	rs12711521	D371Y	$p = 0.04$	MASP2	ENST00000400897	0.86	0.1
1	248084909	rs34508376	M197R	$p=0.01$	OR2T8	ENST00000319968	0.64	0.05
1	248113026	rs10888281	Y289*	—	OR2L8	ENST00000357191	0.94	0.25
1	248129240	rs4478844	V203M	$p=0.00$	OR2AK2	ENST00000366480	0.67	0.05
2	27424636	rs1395	S481F	$p=0.05$	SLC5A6	ENST00000310574	0.74	0.16
5	138720108	rs11242462	W45*	—	SLC23A1	ENST00000508270	0.29	0.80
5	177378959	rs7720935	<i>splice</i>	—	RP11-423H2.3.1	ENST00000507072	0.94	0.40
8	16043667	rs435815	<i>splice</i>	—	MSR1	ENST00000445506	0.11	0.54
19	44932972	rs1434579	G662R	$p=0.04$	ZNF229	ENST00000291187	0.40	0.04
20	2291722	rs6048066	I163L	$p=0.01$	TGM3	ENST00000420960	0.006	0.49

SNPs found in the top 0.2% of XP-SFselect regions, deemed damaging by SIFT (nonsynonymous, with p -value ≤ 0.05) or SnpEff (nonsense or splice-site variant). Frequencies in CEU and YRI populations also shown. Splice site donor mutations are indicated by *splice* in the AA column.