

## SUPPLEMENTARY MATERIAL

### **ProGeM: A framework for the prioritisation of candidate causal genes at molecular quantitative trait loci**

David Stacey<sup>1,\*</sup>, Eric B. Fauman<sup>2</sup>, Daniel Ziemek<sup>3</sup>, Benjamin B. Sun<sup>1</sup>, Eric L. Harshfield<sup>1,4</sup>, Angela M. Wood<sup>1</sup>, Adam S. Butterworth<sup>1</sup>, Karsten Suhre<sup>5</sup>, and Dirk S. Paul<sup>1,\*</sup>

<sup>1</sup> MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK

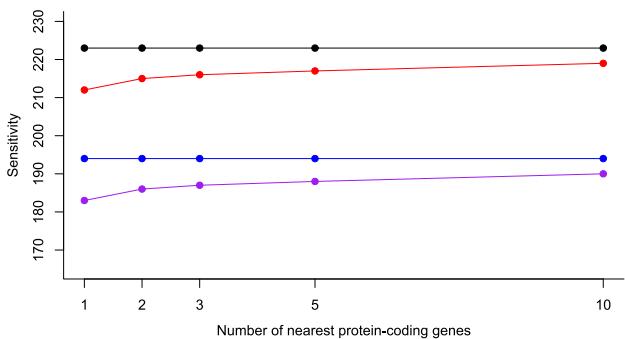
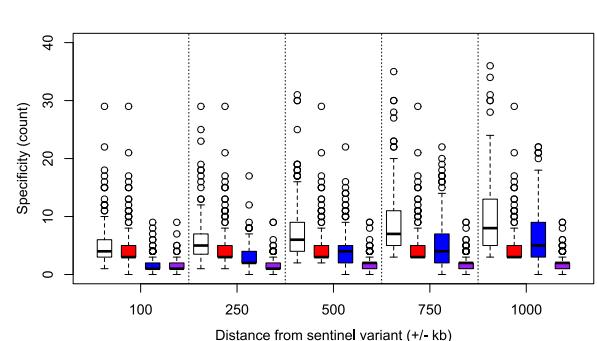
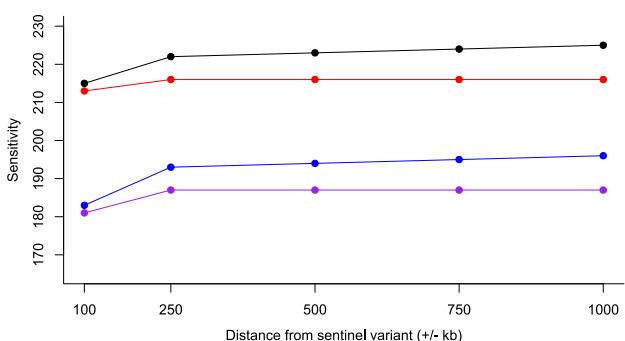
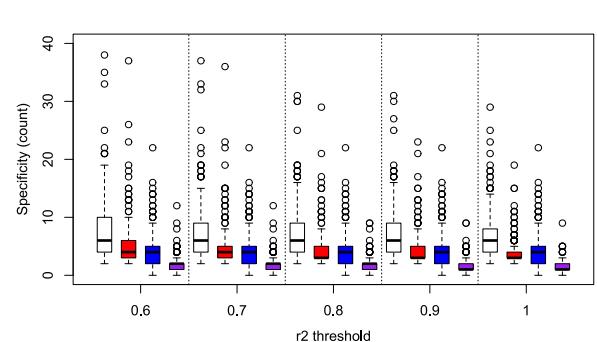
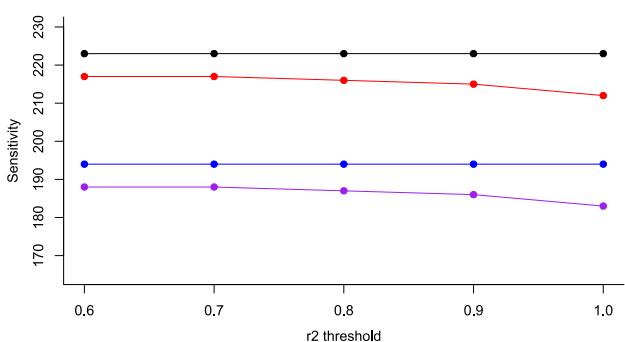
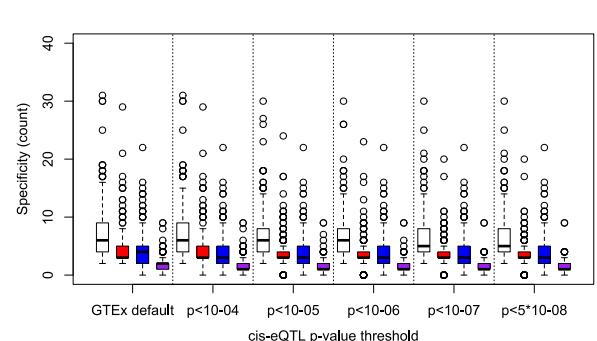
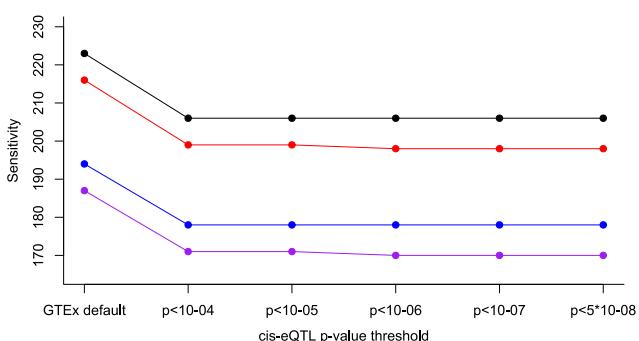
<sup>2</sup> Pfizer Worldwide Research & Development, Genome Sciences & Technologies, Cambridge, MA, USA

<sup>3</sup> Pfizer Worldwide Research & Development, Inflammation & Immunology, Berlin, Germany

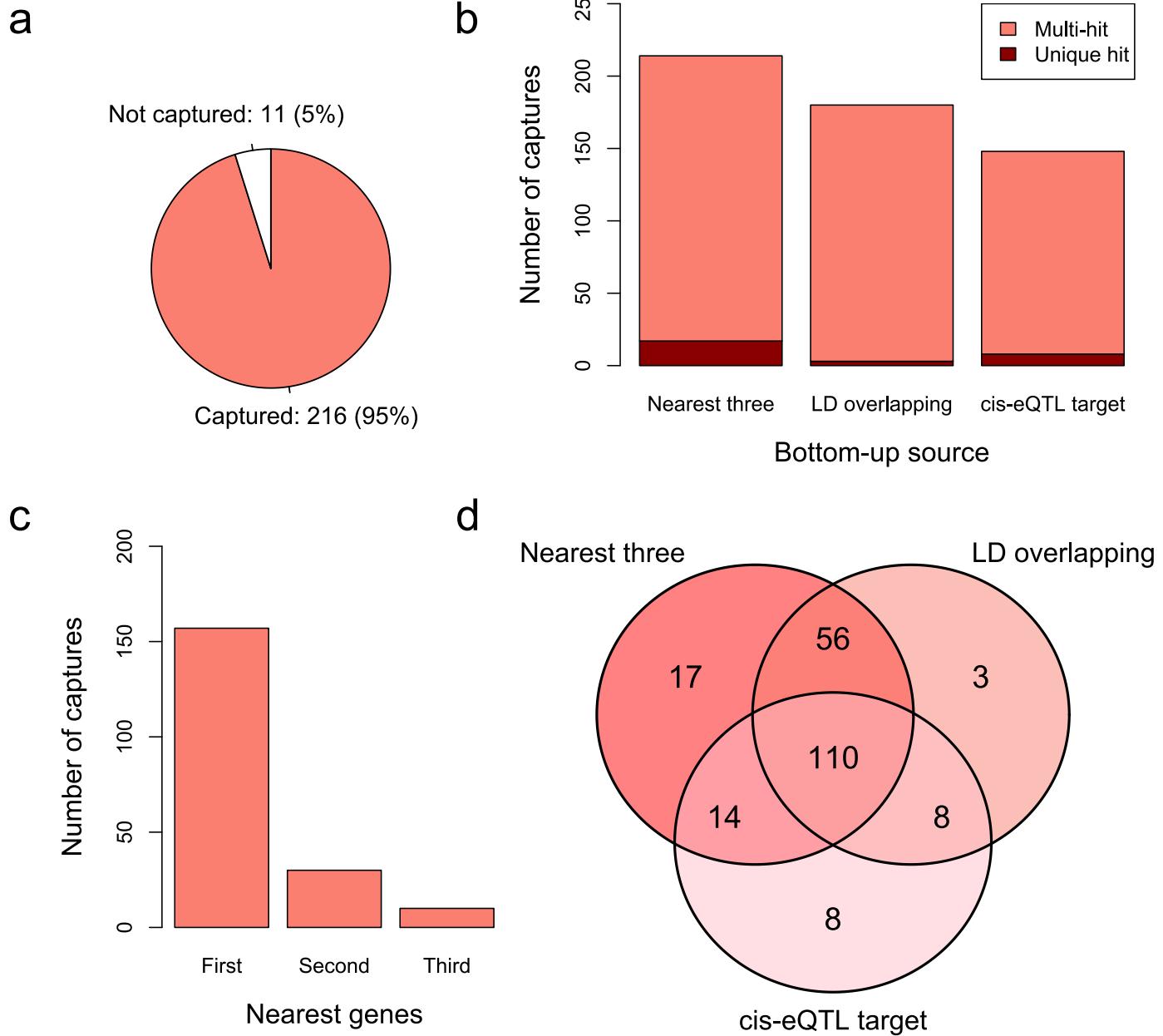
<sup>4</sup> Department of Clinical Neurosciences, University of Cambridge, Cambridge, UK

<sup>5</sup> Department of Physiology and Biophysics, Weill Cornell Medicine-Qatar, Doha, Qatar

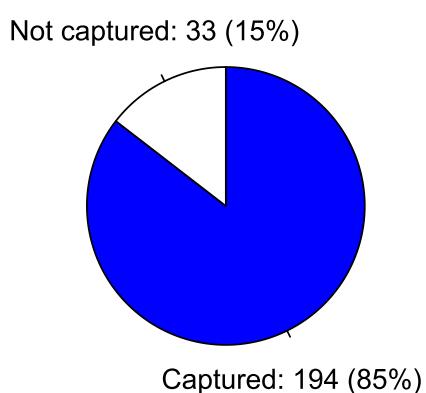
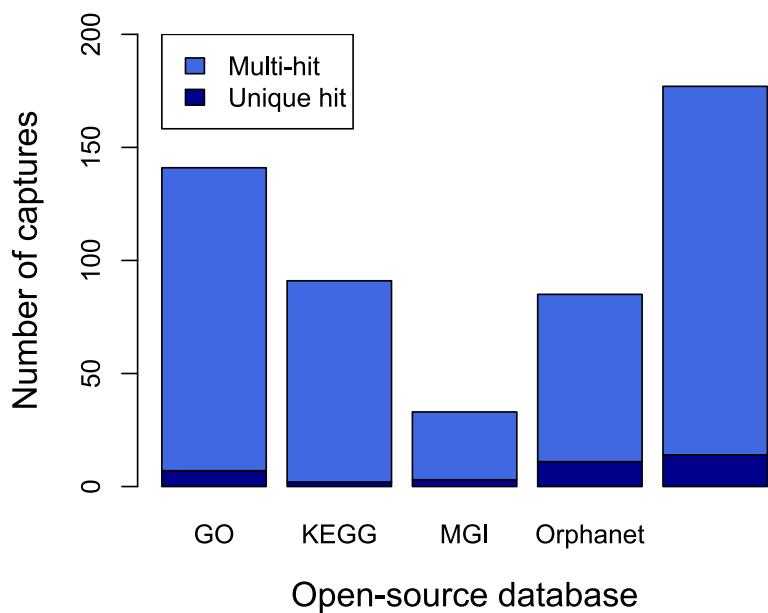
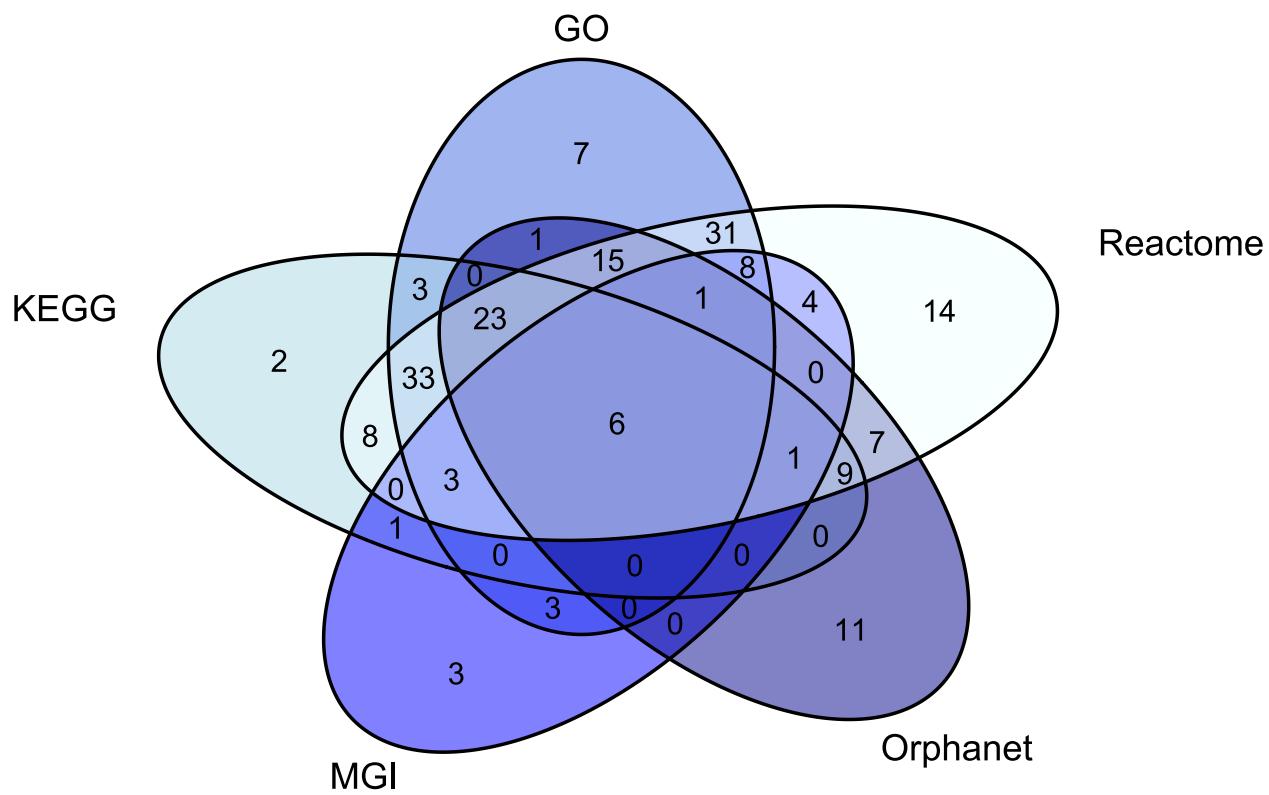
\* To whom correspondence should be addressed. Tel: +44 (0)1223 747217; Email: ds763@medschl.cam.ac.uk.  
Correspondence may also be addressed to. Tel: +44 (0)1223 761918; Email: dsp35@medschl.cam.ac.uk.

**a****b****c****d**

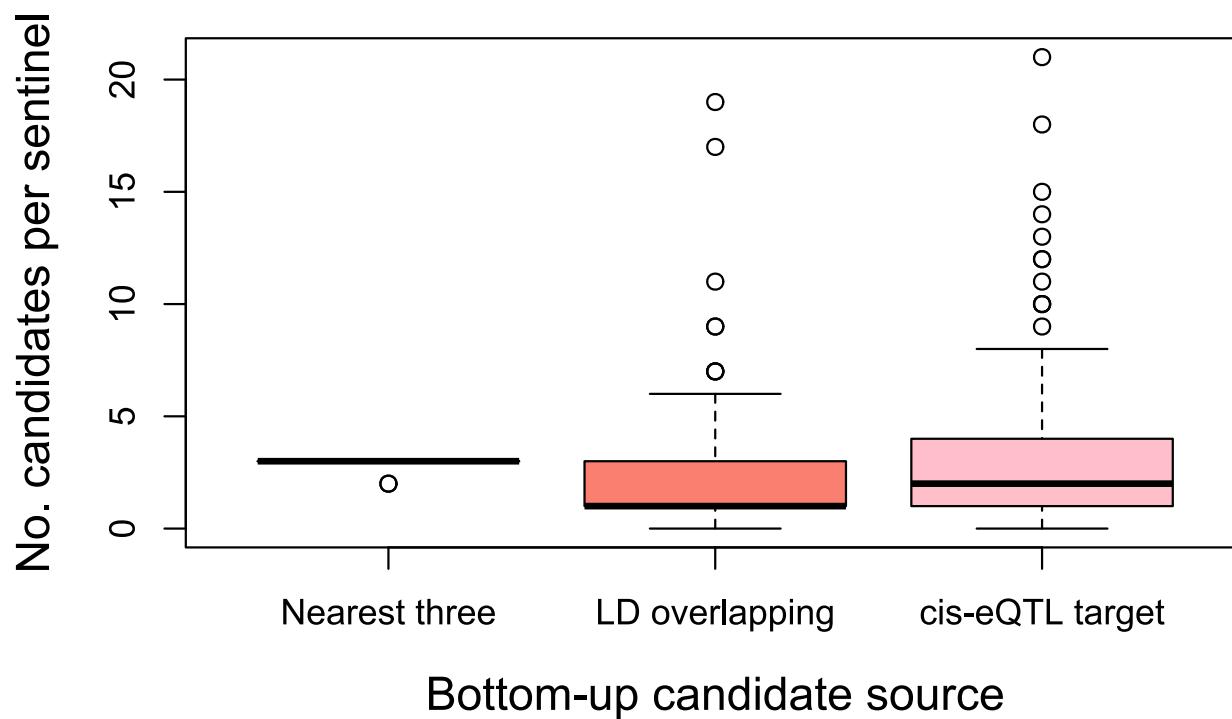
**Supplementary Figure 1. The effects of varying user-defined parameter settings on the sensitivity (line graphs) and specificity (box plots) of ProGeM when applied to a dataset comprising 227 mQTLs.** (a) Number of nearest genes to each sentinel variant (default=3); (b) distance window encompassing each sentinel variant from which to draw candidates (default=500kb); (c) r2 threshold for selecting proxies (default=0.8); and (d) cis-eQTL p-value threshold for selecting cis-eQTL targets as candidate causal genes (default=GTEX default threshold). The box plots show the median and interquartile ranges, with the whiskers extending to 1.5-times the corresponding interquartile range. Data points outside of this range are indicated individually as circles.



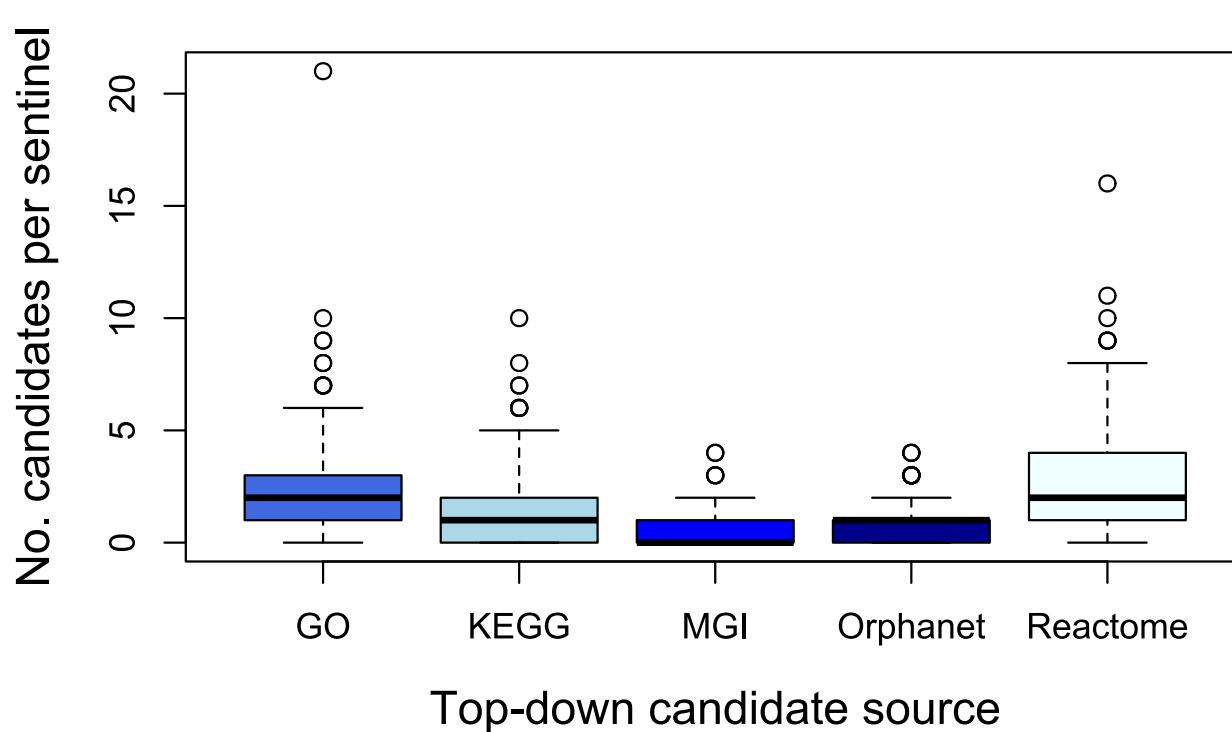
**Supplementary Figure 2. Benchmarking ProGeM: bottom-up sensitivity.** (a) Overall number (percentage) of high-confidence mQTL causal genes captured by the bottom-up component of our framework; (b) Number of high-confidence metabolite QTL causal genes captured by each of the three bottom-up sources utilised by our framework; (c) Number of captured high-confidence causal genes that were one of the three nearest protein-coding genes to their corresponding sentinel variant; (d) Summary of the relative contributions of the three bottom-up sources utilised by our framework towards identifying high-confidence causal genes.

**a****b****c**

**Supplementary Figure 3. Benchmarking ProGeM: top-down sensitivity.** (a) Overall number (percentage) of high-confidence mQTL causal genes captured by the top-down component of ProGeM; (b) and (c) Summaries of the number of high-confidence mQTL causal genes captured either uniquely or concurrently by the five open-source databases utilised by the top-down component of our framework.

**a**

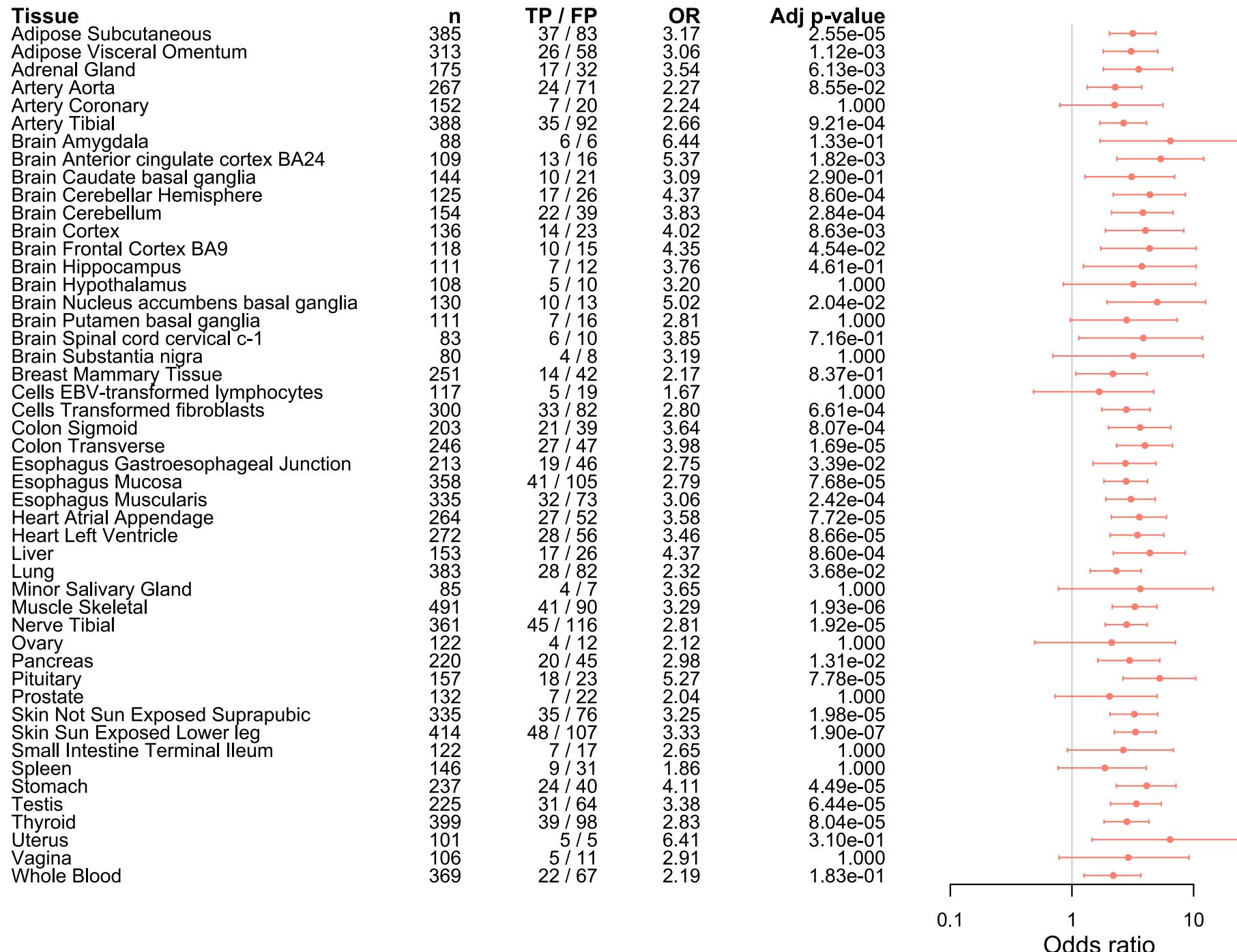
Bottom-up candidate source

**b**

Top-down candidate source

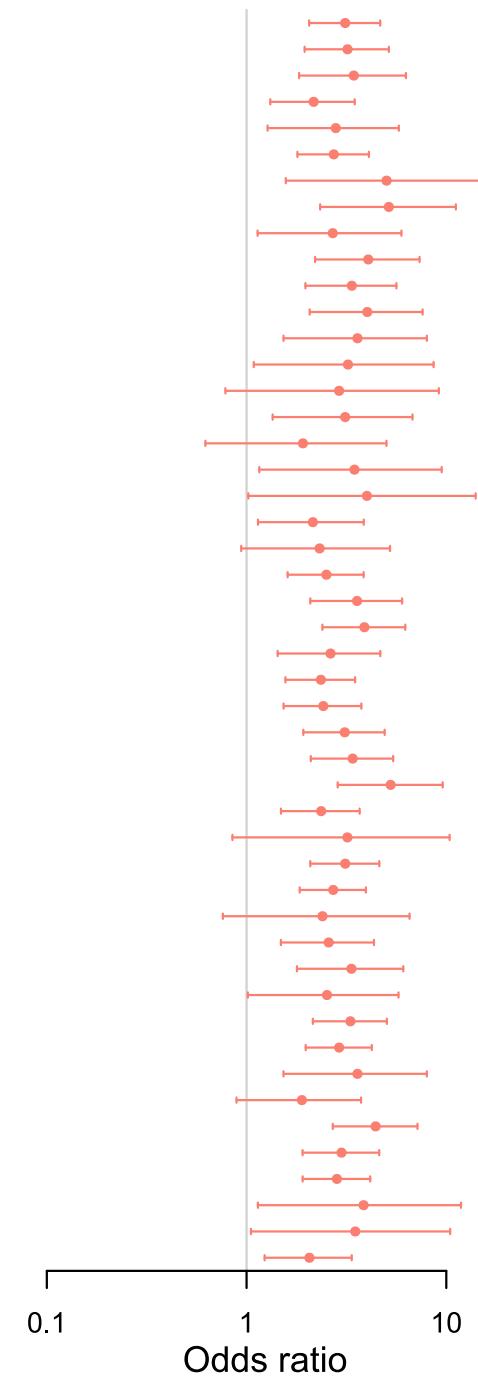
**Supplementary Figure 4. Benchmarking ProGeM: bottom-up and top-down specificity.** Box plots summarising the specificity ("background noise") associated with each of the data sources utilised by the (a) bottom-up and (b) top-down components of ProGeM.

a



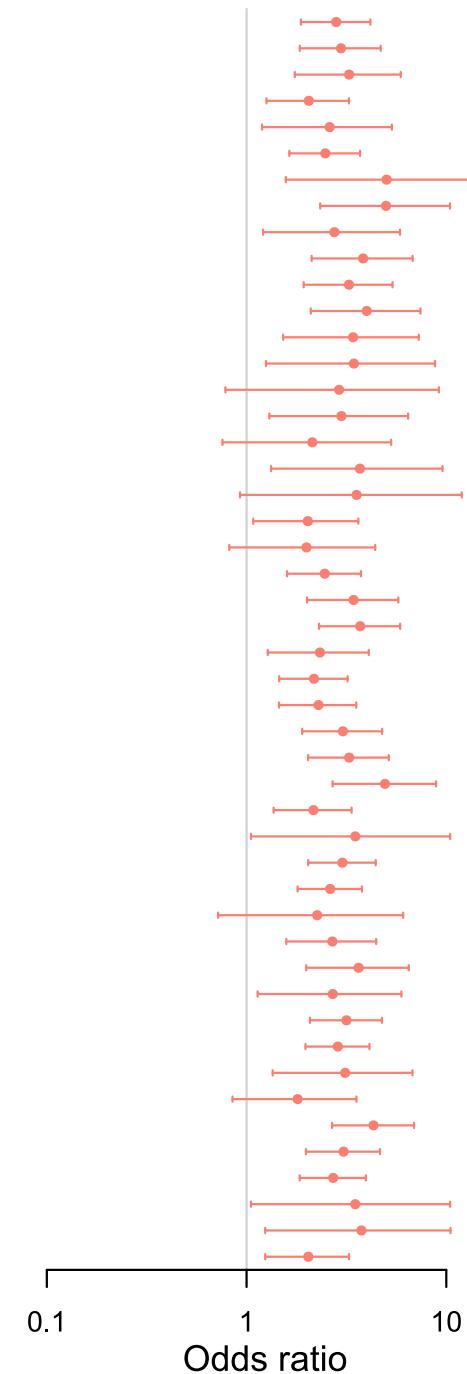
b

Tissue	n	TP / FP	OR	Adj p-value
Adipose Subcutaneous	385	43 / 100	3.12	3.39e-06
Adipose Visceral Omentum	313	30 / 65	3.20	0.000157
Adrenal Gland	175	19 / 37	3.44	0.003433
Artery Aorta	267	27 / 84	2.17	0.073364
Artery Coronary	152	12 / 28	2.80	0.374420
Artery Tibial	388	41 / 107	2.73	9.75e-05
Brain Amygdala	88	7 / 9	5.02	0.150620
Brain Anterior cingulate cortex BA24	109	14 / 18	5.16	0.001265
Brain Caudate basal ganglia	144	10 / 24	2.70	0.894255
Brain Cerebellar Hemisphere	125	21 / 35	4.07	0.000228
Brain Cerebellum	154	26 / 53	3.37	0.000276
Brain Cortex	136	18 / 30	4.02	0.001141
Brain Frontal Cortex BA9	118	11 / 20	3.59	0.081191
Brain Hippocampus	111	7 / 14	3.22	0.835097
Brain Hypothalamus	108	5 / 11	2.91	1.000
Brain Nucleus accumbens basal ganglia	130	11 / 23	3.12	0.187878
Brain Putamen basal ganglia	111	6 / 20	1.92	1.000
Brain Spinal cord cervical c-1	83	7 / 13	3.47	0.628114
Brain Substantia nigra	80	5 / 8	4.00	1.000
Breast Mammary Tissue	251	17 / 52	2.15	0.544134
Cells EBV-transformed lymphocytes	117	9 / 25	2.32	1.000
Cells Transformed fibroblasts	300	35 / 97	2.51	0.002469
Colon Sigmoid	203	26 / 50	3.57	0.000120
Colon Transverse	246	33 / 60	3.89	1.58e-06
Esophagus Gastroesophageal Junction	213	19 / 48	2.63	0.071580
Esophagus Mucosa	358	42 / 126	2.36	0.001303
Esophagus Muscularis	335	33 / 94	2.42	0.005934
Heart Atrial Appendage	264	32 / 72	3.10	0.000118
Heart Left Ventricle	272	32 / 66	3.40	4.07e-05
Liver	153	23 / 30	5.27	3.21e-06
Lung	383	32 / 93	2.36	0.009654
Minor Salivary Gland	85	5 / 10	3.20	1.000
Muscle Skeletal	491	46 / 108	3.12	2.00e-06
Nerve Tibial	361	49 / 132	2.72	2.10e-05
Ovary	122	6 / 16	2.40	1.000
Pancreas	220	23 / 60	2.58	0.022329
Pituitary	157	19 / 38	3.35	0.004526
Prostate	132	9 / 23	2.53	1.000
Skin Not Sun Exposed Suprapubic	335	40 / 87	3.31	2.45e-06
Skin Sun Exposed Lower leg	414	50 / 127	2.91	2.33e-06
Small Intestine Terminal Ileum	122	11 / 20	3.59	0.081191
Spleen	146	12 / 41	1.89	1.000
Stomach	237	33 / 53	4.43	1.56e-07
Testis	225	36 / 85	2.99	6.84e-05
Thyroid	399	47 / 121	2.83	9.17e-06
Uterus	101	6 / 10	3.85	0.716100
Vagina	106	6 / 11	3.50	0.983049
Whole Blood	369	25 / 81	2.06	0.234513



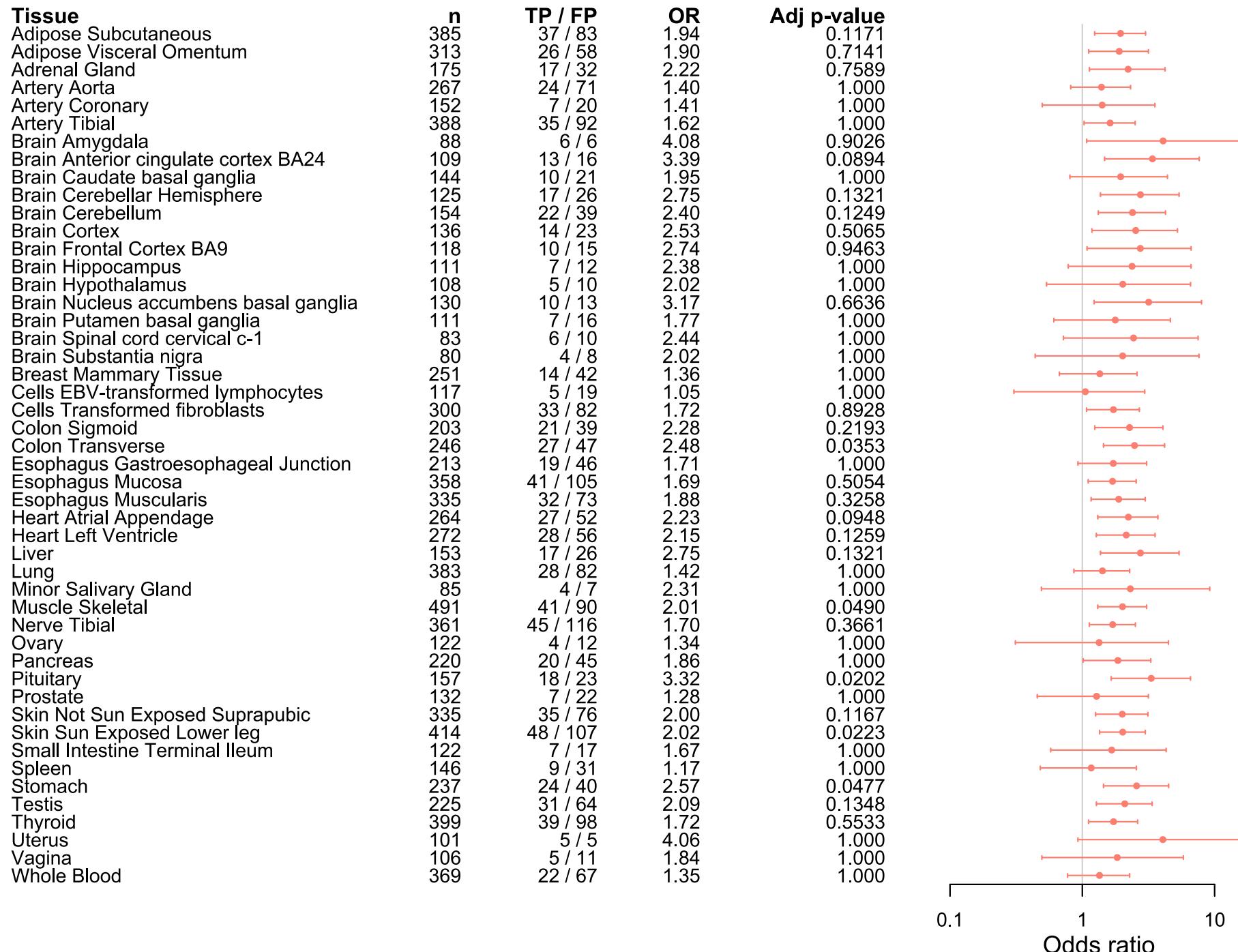
C

Tissue	n	TP / FP	OR	Adj p-value
Adipose Subcutaneous	385	44 / 113	2.81	2.60e-05
Adipose Visceral Omentum	313	32 / 75	2.97	0.000299
Adrenal Gland	175	19 / 39	3.26	0.005916
Artery Aorta	267	28 / 92	2.05	0.160141
Artery Coronary	152	12 / 30	2.61	0.477339
Artery Tibial	388	41 / 117	2.48	0.000676
Brain Amygdala	88	7 / 9	5.02	0.150620
Brain Anterior cingulate cortex BA24	109	15 / 20	4.99	0.000870
Brain Caudate basal ganglia	144	11 / 26	2.75	0.578602
Brain Cerebellar Hemisphere	125	22 / 39	3.83	0.000284
Brain Cerebellum	154	27 / 57	3.26	0.000443
Brain Cortex	136	19 / 32	3.99	0.000746
Brain Frontal Cortex BA9	118	12 / 23	3.42	0.070736
Brain Hippocampus	111	8 / 15	3.45	0.403529
Brain Hypothalamus	108	5 / 11	2.91	1.000
Brain Nucleus accumbens basal ganglia	130	11 / 24	2.98	0.241568
Brain Putamen basal ganglia	111	7 / 21	2.14	1.000
Brain Spinal cord cervical c-1	83	8 / 14	3.70	0.297560
Brain Substantia nigra	80	5 / 9	3.56	1.000
Breast Mammary Tissue	251	17 / 55	2.03	0.997386
Cells EBV-transformed lymphocytes	117	9 / 29	2.00	1.000
Cells Transformed fibroblasts	300	37 / 105	2.46	0.001601
Colon Sigmoid	203	26 / 52	3.43	0.000211
Colon Transverse	246	34 / 65	3.71	2.46e-06
Esophagus Gastroesophageal Junction	213	19 / 54	2.33	0.212216
Esophagus Mucosa	358	43 / 139	2.18	0.006231
Esophagus Muscularis	335	33 / 99	2.29	0.013803
Heart Atrial Appendage	264	33 / 76	3.04	0.000170
Heart Left Ventricle	272	33 / 71	3.26	4.03e-05
Liver	153	23 / 32	4.93	7.43e-06
Lung	383	32 / 101	2.16	0.040382
Minor Salivary Gland	85	6 / 11	3.50	0.983049
Muscle Skeletal	491	48 / 117	3.02	1.60e-06
Nerve Tibial	361	52 / 146	2.62	1.97e-05
Ovary	122	6 / 17	2.26	1.000
Pancreas	220	25 / 63	2.69	0.008274
Pituitary	157	21 / 39	3.64	0.000807
Prostate	132	10 / 24	2.70	0.894255
Skin Not Sun Exposed Suprapubic	335	42 / 96	3.16	3.74e-06
Skin Sun Exposed Lower leg	414	54 / 141	2.86	1.26e-06
Small Intestine Terminal Ileum	122	11 / 23	3.12	0.187878
Spleen	146	12 / 43	1.80	1.000
Stomach	237	35 / 58	4.32	9.00e-08
Testis	225	39 / 91	3.06	1.76e-05
Thyroid	399	49 / 132	2.72	2.10e-05
Uterus	101	6 / 11	3.50	0.983049
Vagina	106	7 / 12	3.76	0.460760
Whole Blood	369	27 / 89	2.04	0.150785



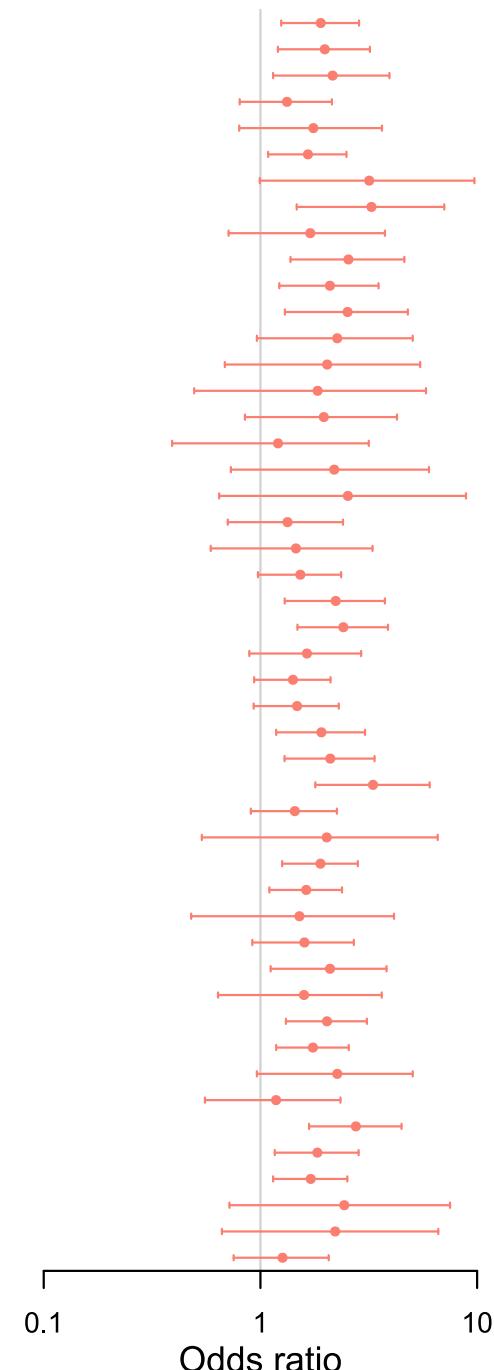
**Supplementary Figure 5. Comparison of tissue-specific cis-eQTL data from GTEx v7 (n=48 tissues) for distinguishing true from false positive causal genes at 227 independent mQTLs.** Odds ratios and 95% confidence intervals are indicated for cis-eQTL targets of (a) sentinel variants, (b) proxy variants, and (c) either sentinel or proxy variants. The background gene set comprised of all candidate causal genes (both top-down and bottom-up) highlighted by ProGeM. Fisher's exact test was used throughout, and Bonferroni corrected (48 tests) *p*-values are indicated. The number of true positive (TP) and false positive (FP) causal genes identified by each characteristic, as well as the number of samples for each tissue assayed by the GTEx consortium, are also indicated.

a



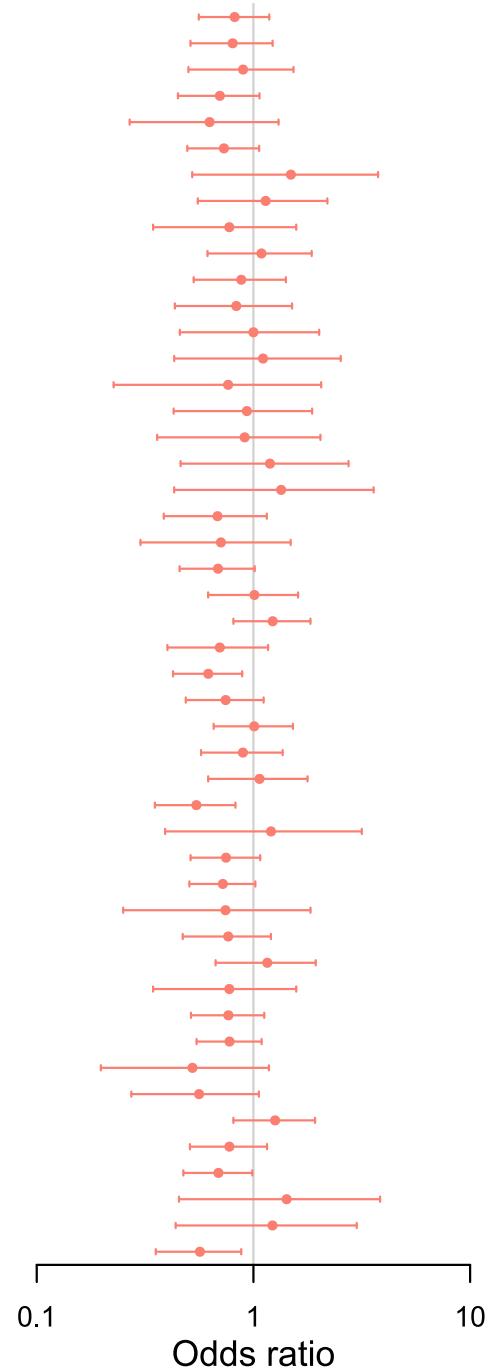
b

Tissue	n	TP / FP	OR	Adj p-value
Adipose Subcutaneous	385	43 / 100	1.90	0.10905
Adipose Visceral Omentum	313	30 / 65	1.98	0.22337
Adrenal Gland	175	19 / 37	2.15	0.71501
Artery Aorta	267	27 / 84	1.33	1.000
Artery Coronary	152	12 / 28	1.76	1.000
Artery Tibial	388	41 / 107	1.66	0.70801
Brain Amygdala	88	7 / 9	3.18	1.000
Brain Anterior cingulate cortex BA24	109	14 / 18	3.25	0.11344
Brain Caudate basal ganglia	144	10 / 24	1.70	1.000
Brain Cerebellar Hemisphere	125	21 / 35	2.55	0.08223
Brain Cerebellum	154	26 / 53	2.09	0.24251
Brain Cortex	136	18 / 30	2.52	0.22375
Brain Frontal Cortex BA9	118	11 / 20	2.26	1.000
Brain Hippocampus	111	7 / 14	2.03	1.000
Brain Hypothalamus	108	5 / 11	1.84	1.000
Brain Nucleus accumbens basal ganglia	130	11 / 23	1.96	1.000
Brain Putamen basal ganglia	111	6 / 20	1.21	1.000
Brain Spinal cord cervical c-1	83	7 / 13	2.19	1.000
Brain Substantia nigra	80	5 / 8	2.53	1.000
Breast Mammary Tissue	251	17 / 52	1.33	1.000
Cells EBV-transformed lymphocytes	117	9 / 25	1.46	1.000
Cells Transformed fibroblasts	300	35 / 97	1.53	1.000
Colon Sigmoid	203	26 / 50	2.23	0.12502
Colon Transverse	246	33 / 60	2.41	0.01623
Esophagus Gastroesophageal Junction	213	19 / 48	1.64	1.000
Esophagus Mucosa	358	42 / 126	1.41	1.000
Esophagus Muscularis	335	33 / 94	1.48	1.000
Heart Atrial Appendage	264	32 / 72	1.91	0.30545
Heart Left Ventricle	272	32 / 66	2.10	0.09970
Liver	153	23 / 30	3.31	0.00307
Lung	383	32 / 93	1.44	1.000
Minor Salivary Gland	85	5 / 10	2.02	1.000
Muscle Skeletal	491	46 / 108	1.89	0.07307
Nerve Tibial	361	49 / 132	1.63	0.52324
Ovary	122	6 / 16	1.51	1.000
Pancreas	220	23 / 60	1.60	1.000
Pituitary	157	19 / 38	2.09	0.76557
Prostate	132	9 / 23	1.59	1.000
Skin Not Sun Exposed Suprapubic	335	40 / 87	2.03	0.06190
Skin Sun Exposed Lower leg	414	50 / 127	1.75	0.18626
Small Intestine Terminal Ileum	122	11 / 20	2.26	1.000
Spleen	146	12 / 41	1.18	1.000
Stomach	237	33 / 53	2.76	0.00166
Testis	225	36 / 85	1.83	0.36262
Thyroid	399	47 / 121	1.71	0.30238
Uterus	101	6 / 10	2.44	1.000
Vagina	106	6 / 11	2.21	1.000
Whole Blood	369	25 / 81	1.26	1.000



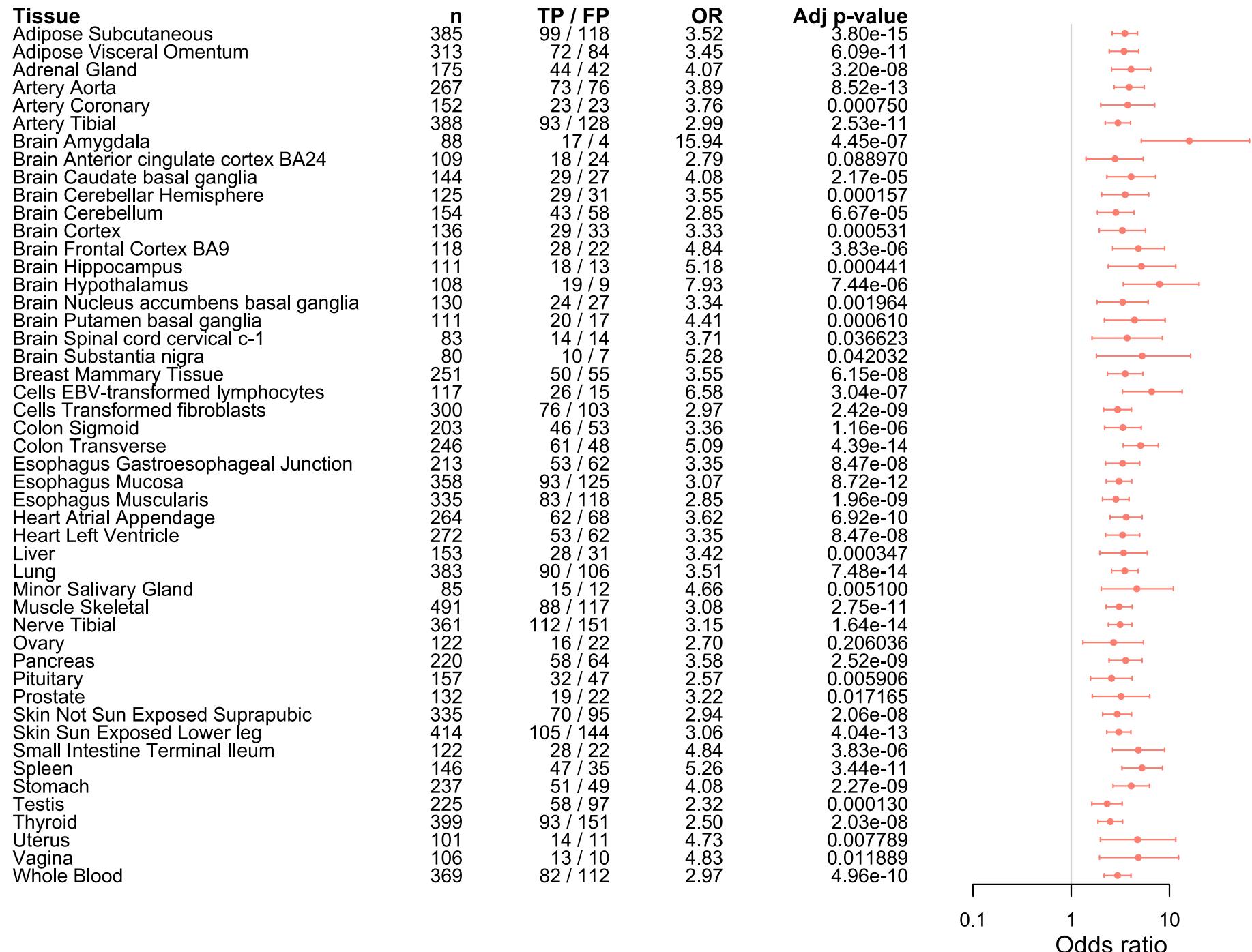
C

Tissue	n	TP / FP	OR	Adj p-value
Adipose Subcutaneous	385	47 / 219	0.82	1.000
Adipose Visceral Omentum	313	32 / 154	0.80	1.000
Adrenal Gland	175	19 / 84	0.90	1.000
Artery Aorta	267	32 / 172	0.70	1.000
Artery Coronary	152	9 / 56	0.63	1.000
Artery Tibial	388	44 / 224	0.73	1.000
Brain Amygdala	88	7 / 19	1.49	1.000
Brain Anterior cingulate cortex BA24	109	13 / 46	1.14	1.000
Brain Caudate basal ganglia	144	10 / 51	0.77	1.000
Brain Cerebellar Hemisphere	125	20 / 74	1.09	1.000
Brain Cerebellum	154	25 / 112	0.88	1.000
Brain Cortex	136	15 / 71	0.83	1.000
Brain Frontal Cortex BA9	118	11 / 44	1.00	1.000
Brain Hippocampus	111	8 / 29	1.11	1.000
Brain Hypothalamus	108	5 / 26	0.76	1.000
Brain Nucleus accumbens basal ganglia	130	11 / 47	0.93	1.000
Brain Putamen basal ganglia	111	8 / 35	0.91	1.000
Brain Spinal cord cervical c-1	83	8 / 27	1.19	1.000
Brain Substantia nigra	80	6 / 18	1.34	1.000
Breast Mammary Tissue	251	19 / 107	0.68	1.000
Cells EBV-transformed lymphocytes	117	9 / 50	0.71	1.000
Cells Transformed fibroblasts	300	39 / 210	0.69	1.000
Colon Sigmoid	203	27 / 107	1.01	1.000
Colon Transverse	246	40 / 135	1.23	1.000
Esophagus Gastroesophageal Junction	213	20 / 110	0.70	1.000
Esophagus Mucosa	358	48 / 273	0.62	0.366
Esophagus Muscularis	335	36 / 183	0.74	1.000
Heart Atrial Appendage	264	36 / 143	1.01	1.000
Heart Left Ventricle	272	33 / 145	0.89	1.000
Liver	153	22 / 83	1.07	1.000
Lung	383	32 / 209	0.55	0.123
Minor Salivary Gland	85	6 / 20	1.21	1.000
Muscle Skeletal	491	48 / 239	0.75	1.000
Nerve Tibial	361	56 / 282	0.72	1.000
Ovary	122	6 / 32	0.74	1.000
Pancreas	220	27 / 136	0.76	1.000
Pituitary	157	22 / 77	1.16	1.000
Prostate	132	10 / 51	0.77	1.000
Skin Not Sun Exposed Suprapubic	335	42 / 207	0.77	1.000
Skin Sun Exposed Lower leg	414	59 / 282	0.78	1.000
Small Intestine Terminal Ileum	122	7 / 52	0.52	1.000
Spleen	146	12 / 82	0.56	1.000
Stomach	237	35 / 115	1.26	1.000
Testis	225	37 / 182	0.77	1.000
Thyroid	399	49 / 258	0.69	1.000
Uterus	101	6 / 17	1.42	1.000
Vagina	106	7 / 23	1.22	1.000
Whole Blood	369	28 / 180	0.57	0.432

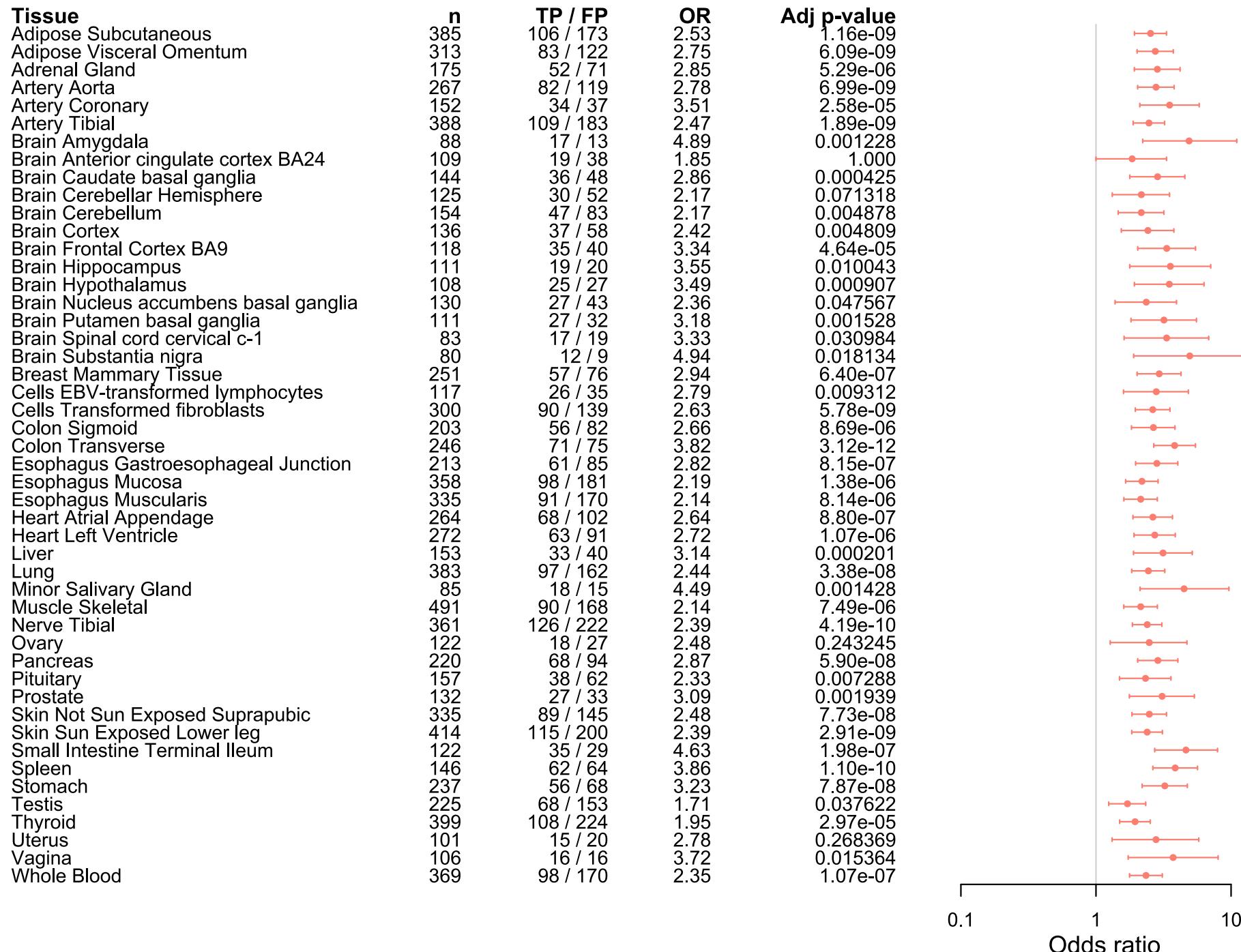


**Supplementary Figure 6. Comparison of tissue-specific cis-eQTL data from GTEx v7 (n=48 tissues) for distinguishing true from false positive causal genes at 227 independent mQTLs using the bottom-up background gene set.** Odds ratios and 95% confidence intervals are indicated for cis-eQTL targets of (a) sentinel variants, (b) proxy variants, and (c) either sentinel or proxy variants are indicated. The background gene set comprised of all bottom-up candidate causal genes highlighted by ProGeM. Fisher's exact test was used throughout, and Bonferroni corrected (48 tests) *p*-values are indicated. The number of true positive (TP) and false positive (FP) causal genes identified by each characteristic, as well as the number of samples for each tissue assayed by the GTEx consortium, are also indicated.

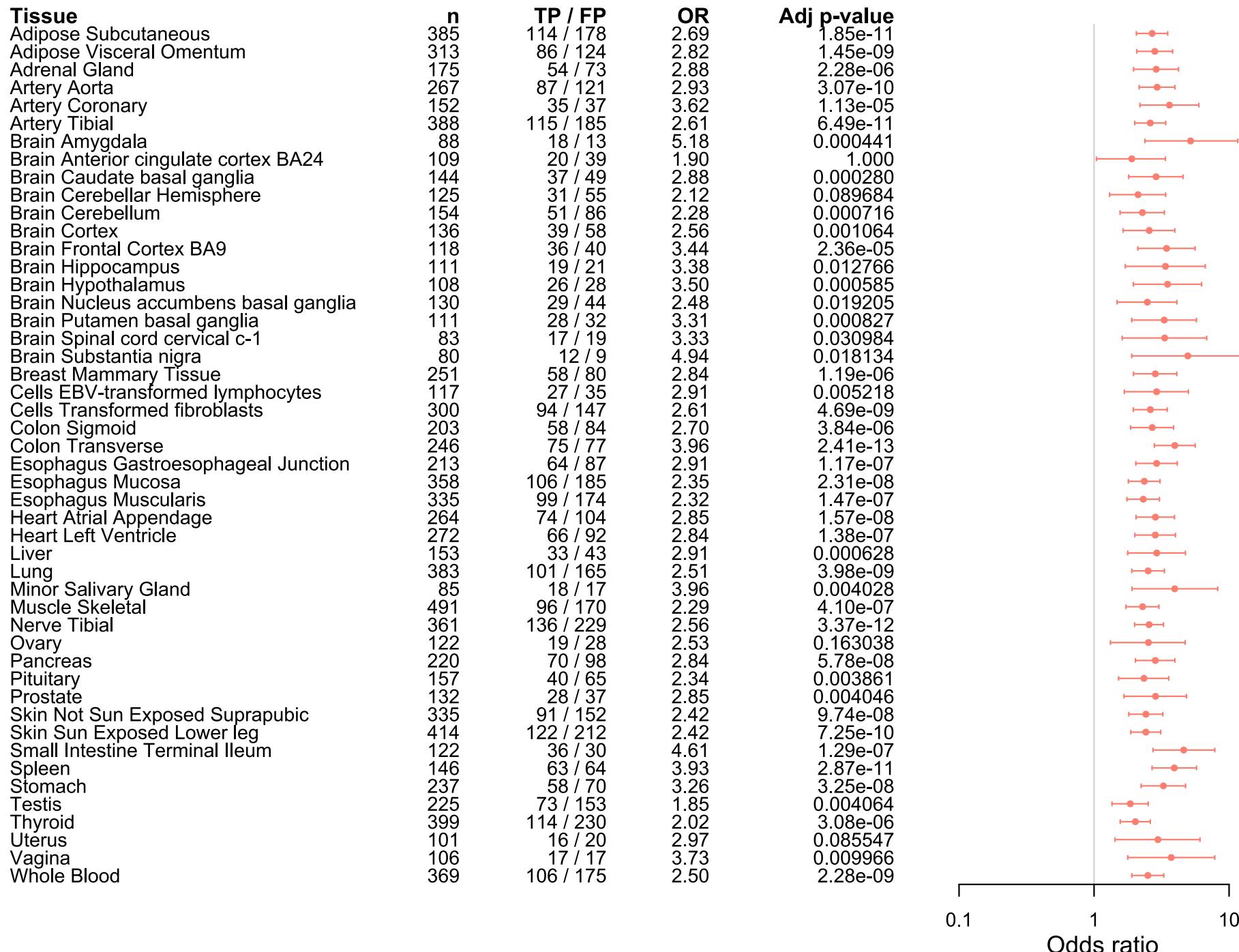
a



b



C



**Supplementary Figure 7. Comparison of tissue-specific cis-eQTL from GTEx v7 data (n=48 tissues) for distinguishing true from false positive causal genes at 562 independent cis-pQTLs.** Odds ratios and 95% confidence intervals are indicated for cis-eQTL targets of (a) sentinel variants, (b) proxy variants, and (c) either sentinel or proxy variants are indicated. The background gene set comprised of all bottom-up candidate causal genes highlighted by ProGeM. Fisher's exact test was used throughout, and Bonferroni corrected (48 tests) *p*-values are indicated. The number of true positive (TP) and false positive (FP) causal genes identified by each characteristic, as well as the number of samples for each tissue assayed by the GTEx consortium, are also indicated.