# Tissue-specific genes as an underutilized resource in drug discovery

Maria Ryaboshapkina MSc[1,*], Mårten Hammar PhD[1]

[1]Translational Science, Cardiovascular, Renal and Metabolism, IMED Biotech Unit,

AstraZeneca, Gothenburg, Sweden

*maria.ryaboshapkina@astrazeneca.com


AstraZeneca

Pepparedsleden 1

431 50 Mölndal

Sweden

# ABSTRACT

Tissue-specific genes are believed to be good drug targets due to improved safety. Here we show that this intuitive notion is not reflected in phase 1 and 2 clinical trials, despite the historic success of tissue-specific targets and their 2.3-fold overrepresentation among targets of marketed non-oncology drugs. We compare properties of tissue-specific genes and drug targets. We show that tissue-specificity of the target may also be related to efficacy of the drug. The relationship may be indirect (enrichment in Mendelian disease genes) or direct (elevated ability to spread perturbations in human protein-protein interactome for tissue-specifically produced enzymes and secreted proteins). Reduced evolutionary conservation of tissue-specific genes may represent a bottleneck for drug projects, prompting development of novel models with smaller evolutionary gap to humans. We highlight numerous open opportunities to use tissue-specific genes in drug research and hope that the current study will facilitate discovery efforts.

50    Narrow expression in one or a few tissues is considered desirable for drug targets due

51    to reduced risk of side effects[1,2]. Genes with narrow expression are often called 'tissue-

52    specific' or 'tissue-enriched'. Studies on microarray[3-5] and a combination of RNA-

53    sequencing and proteomics data[6,7] confirm that targets of marketed drugs are biased

54    towards tissue-specific genes. To the best of our knowledge, the first quantitative

55    estimate was published in 2008. Dezso et al. demonstrated that tissue-specific genes

56    are twice more likely to become drug targets than broadly expressed house-keeping

57    genes[8]. Yang et al. confirmed a 1.7-fold higher likelihood in 2016[9]. Dezso et al.

58    observed that tissue-specific genes may represent attractive drug targets due to their

59    role in tissue biology and disease (e.g., brain-specific *GABRB2*, a receptor for the

60    inhibitory neuromediator gamma-aminobutyric acid*,* is a target of sedative agents)[8].

61    These studies assessed tissue-specificity in healthy tissues. Their findings also

62    extrapolate to diseased tissues because targets of marketed and phase 3 drugs are

63    expressed in disease-relevant tissues even in the healthy state in 87% of the cases[10].

64    Also, substantial efforts are dedicated to cataloguing tissue-specific genes such as

65    databases TiGER (2008)[11], TiSGeD (2010)[12], VeryGene (2011)[13] and TissGDB

66    (2018)[14]. Thus, systematic studies showing a significant overrepresentation of tissue-

67    specific genes among drug targets and comprehensive resources have been available

68    since 2008. The average time from a lead compound to entering phase 1 clinical trial

69    is 31.2 months[15]. Let's assume that validation of the biological hypothesis and

70    identification of a lead compound take an equally long time. Then, ten years are

71    sufficient for the findings of basic research to find reflection in early phase clinical trials.

72    Now is good time to test if the industry took advantage of omics studies and pursued

73    tissue-specific targets or not.

74    In this study, we examine the prevalence of tissue-specific genes among targets of

75    marketed drugs and drugs in clinical trials. We also investigate properties of tissue-

76    specific genes compared to drug targets. Why are these questions important to

77    address? If tissue-specific targets are not actively pursued in early clinical trials, such

78    study would raise awareness of open opportunities. Opportunities to discover new

79    targets are not exhausted. A recent study by Oprea et al. indicates that only 3% of

80    human proteins are targeted by marketed or clinical trial drugs ("Tclin") whereas 35%

81    have an unknown biological function and are not actively studied ("Tdark")[16]. Also, an

82    important parallel exists between tissue-specific genes and targets of marketed drugs.

83    As first demonstrated in 2004, tissue-specific genes are enriched in Mendelian

84    disorder genes[17]. The enrichment was confirmed by Yang et al. in 2016[9]. 53% targets

85    of marketed drugs are implicated in Mendelian disorders[18]. Drugs targeting genes with

86    a genetic link to human disease are less likely to fail in clinical trials due to lack of

87    efficacy[19]. Thus, there may be a relationship between tissue-specificity of the target

88    and efficacy of the drug. In fact, a recent study by Rouillard, Hurle and Agarwal

89    concentrated on identification of omics features distinguishing targets that succeeded

90    and failed in phase 3 trials for non-oncology diseases[20]. Phase 3 trial failures were

91    enriched in failures due to lack of efficacy. Rouillard and colleagues limited their

92    analysis to drugs with a single mechanism-of-action target and demonstrated that

93    narrow expression profile of a drug target is a robust predictor of success in phase 3[20].

94    If we understand the relationship between tissue-specificity and efficacy and apply this

95    knowledge to identify new, and not only tissue-specific, targets, we may reduce

96    attrition rates in the clinic.

97    Here we confirm a 1.8-fold overrepresentation of tissue-specific genes among targets

98    of marketed drugs compared to all protein-coding genes. The enrichment is 2.3-fold

99   when non-oncology drug targets are considered separately. We observe that this

100  historic success of tissue-specific targets is not reflected in early clinical trials neither

101  for oncology nor for non-oncology diseases. We find two factors, that could be related

102  to efficacy of drugs targeting tissue-specific genes. First, we confirm enrichment in

103  disease genes among tissue-specific genes. Second, we find that tissue-specific

104  enzymes and secreted proteins have higher ability to spread perturbations in

105  topological analysis of human protein-protein interactome. The limiting factor for

106  development of tissue-specific targets may be the reduced conservation of tissue-

107  specific genes between humans and murine models and the associated challenges in

108  preclinical research. We conclude that tissue-specific genes are a promising source

109  for target discovery and that the translational challenges may be circumvented through

110  creation of humanized models.

## RESULTS

111

112  Our results section is structured as follows. We investigate the prevalence of tissue-

113  specific genes among targets of candidate and marketed drugs. Next, we explore

114  properties that may explain depletion of tissue-specific genes among targets of drugs

115  in early clinical trials and their overrepresentation among targets of marketed drugs.

116  Finally, we highlight open opportunities to develop tissue-specific genes as drug

117  targets.

118  We talk about genes as drug targets because the previous studies demonstrated

119  enrichment in tissue-specific genes among drug targets based on mRNA

120  expression[8,9]. We also define tissue-specificity based on RNA-sequencing data. We

121  assume that the messenger RNAs are translated to their protein products, which, in

122  turn, interact with the drugs. The concordance between gene expression and protein

5

123  abundance is debated[21,22], but a recent Ribo-seq study in rat suggests that 70 (heart)

124  to 85% (liver) of transcribed mRNA are forwarded to translation[23].

## Prevalence of tissue-specific drug targets

126  We applied peak-based definitions of tissue-specificity. We computed per-tissue Z-

127  scores for each gene and defined tissue-specificity at nine increasingly stringent

128  constraints: $Z_{\text{second largest}} < 1/x * Z_{\text{max}}$, where $Z_{max}$ denoted the Z-score in the tissue with

129  the highest expression, $Z_{\text{second largest}}$ denoted the Z-score in tissue with the second

130  highest expression and x was an integer from 2 to 10 (**Fig. 1)**. Such definitions allowed

131  genes to be expressed in multiple tissues, as long as there was a clear "peak" in the

132  tissue with the highest expression compared to all other tissues.

133  Tissue-specific genes constituted a small fraction of all human protein-coding genes

134  (**Supplementary Data 1**). The most liberal definition x = 2 resulted in 4,573 of 18,377,

135  24.9% tissue-specific genes, while only 557 of 18,377, 3.0% genes satisfied the most

136  stringent definition x = 10. If tissue-specificity was irrelevant for drug target discovery,

137  the proportions of tissue-specific genes among drug targets would follow the

138  'background' distribution among all protein-coding genes. By contrast, we observed

139  increasingly stronger deviations from the 'background' distribution with increasingly

140  stringent definitions of tissue-specificity (**Fig. 2**). Targets of phase 1 drugs were

141  significantly depleted of tissue-specific genes even at the liberal x = 2 (54 of 331,

142  16.3% < 4,573 of 18,377, 24.9%, Fisher test, p-value $1*10^{-4}$). Proportions of tissue-

143  specific genes among targets of phase 2 drugs followed the 'background' distribution

144  among all protein-coding genes. By contrast, targets of phase 3 drugs and marketed

145  drugs were significantly enriched in tissue-specific genes starting from x = 6 (phase 3:

146  39 of 410, 9.5% > 1,018 of 18,377, 5.5%, 1.7-fold enrichment, p-value 0.001;

147  marketed: 70 of 691, 10.1% > 1,018 of 18,377, 5.5%, 1.8-fold enrichment, p-value

148    $2*10^{-6}$). Targets of withdrawn drugs were also enriched in tissue-specific genes at x =

149    8 to 10. The reason for withdrawal from the market was toxicity with few exceptions

150    like unintended use for self-poisoning (barbiturates) and lack of efficacy (drotrecogin

151    alpha). Targets of withdrawn drugs had 95% overlap (57 of 60) with targets of

152    marketed drugs. Hence, withdrawal of these drugs from the market could not be

153    uniquely attributed to their mechanism-of-action targets. For example, cholinergic

154    nicotinic receptors *CHRNA1*, *CHRND* and *CHRNG* are targets of curare-like

155    neuromuscular blocking agents. Rapacuronium bromide was withdrawn from the

156    market due to adverse events while other drugs like vecuronium continue to be used.

157    We used x = 6 to define tissue-specificity in all subsequent analyses, because the

158    enrichment in tissue-specific genes among targets of marketed drugs became

159    significant at this constraint.

160    The overlap between targets of withdrawn and marketed drugs motivated us to

161    examine 'recycling' of drug targets. Target genes can re-enter clinical trials when new

162    drugs are developed for the same (e.g., generations of $H_2$ histamine receptor *HRH2*

163    blockers as anti-ulcer drugs) or a novel indication. For example, *IGF1R* is targeted by

164    recombinant insulin growth factor 1 Mecasermin for growth failure in children

165    (marketed agonist drug) and is evaluated as target for treatment of solid tumours

166    (antagonist drug PL-225B in phase 1 trial NCT01779336). 431 of all 691 targets,

167    62.4% (**Fig. 3a**) and only 24 of 70 tissue-specific targets, 34.3% (**Fig. 3b**) were reused

168    in clinical trials. Thus, tissue-specific genes were 'recycled' 1.8 times less frequently

169    than drug targets overall (Fisher test p-value $5.6*10^{-6}$). Furthermore, tissue-specific

170    genes represented an older subset of drug targets (**Fig. 3c**), although the difference

171    was not statistically significant.

7

172    In summary, tissue-specific genes, satisfying x = 6, were 1.8 times more likely to

173    become targets of marketed drug than all protein-coding genes. However, they were

174    not actively explored in phase 1 and 2 clinical trials. We investigated possible

175    explanations for these trends.

## Disease indication as a confounding factor

177    Targets for oncology drugs are selected following different paradigms than targets for

178    non-oncology drugs. For example, traditional cytotoxic agents aim to induce cell death

179    or inhibit growths through core cell processes, that are carried out by ubiquitously

180    expressed targets like DNA topoisomerase II (etoposide). Oncology drugs have

181    different safety profiles from non-oncology drugs, with more side effects being

182    tolerated. Also, some drugs target cancer-specific mutant proteins, which are not

183    captured by gene expression analysis on healthy tissue. For example, vemurafenib

184    targets    mutated    *BRAF*    in    melanoma    according    to    the    FDA    label

185    (www.accessdata.fda.gov/drugsatfda_docs/label/2017/202429s012lbl.pdf).

186    Targets of phase 1 drugs were predominantly investigated for oncology indications:

187    253 of 331, 76,4%. Targets of phase 2 drugs displayed an almost balanced

188    representation of targets for oncology - 239 of 553, 43.2% - and non-oncology

189    indications - 314 of 553, 56.8%. By contrast, most targets of phase 3 drugs - 265 of

190    410, 64.6% - and of marketed drugs – 479 of 691, 69.3% - were developed for non-

191    oncology indications. Prevalence of tissue-specific genes among targets of clinical trial

192    and marketed drugs was confounded by disease indications.

193    Hence, we examined oncology and non-oncology targets separately (**Supplementary**

194    **Fig. 1**). Targets of phase 1 drugs were depleted of tissue-specific genes irrespective

195    of disease indication. The discrepancies between oncology and non-oncology targets

196    started to emerge in phase 2. Targets of marketed non-oncology drugs displayed a

8

197 2.3-fold overrepresentation in tissue-specific genes (at x = 6: 61 of 479, 12.7% > 1,018

198 of 18,377, 5.5%, Fisher test, p-value 3.6*$10^{-9}$), which was stronger compared to pooled

199 analysis for all disease indications. By contrast, tissue-specific genes were

200 underrepresented among targets of oncology drugs.

## Insights from evolutionary biology and population genetics

202 Evolutionary properties may explain the underrepresentation of tissue-specific targets

203 in early clinical trials. Wenhua Lv et al. demonstrated that targets of FDA-approved

204 drugs are more evolutionary conserved than non-target genes[24]. By contrast, **i**n 2004,

205 Winter, Goodstasdt and Ponting investigated expression of 4,960 human genes in 27

206 tissues and demonstrated that tissue-specific genes are less evolutionary conserved

207 than broadly expressed genes using $K_a/K_s$ ratios[17]. To clarify, $K_a/K_s$ is the rate of

208 nonsynonymous $K_a$ to synonymous $K_s$ amino acid changes in a pair of orthologs. Low

209 $K_a/K_s$ implies that nonsynonymous changes are selected out, while $K_a/K_s$ exceeding 1

210 may indicate that changes are favored and retained as in immune genes adapting to

211 new pathogens[25].

212 We revisited the analysis with the current larger data set. We examined $K_a/K_s$ ratios

213 for human protein-coding genes and their mouse counterparts because mice are the

214 most common species in preclinical research. We confirmed opposite patterns for

215 $K_a/K_s$ ratios of tissue-specific genes and drug targets (**Fig. 4a**). Tissue-specific genes

216 had significantly higher $K_a/K_s$ than all protein-coding genes (Mann-Whitney U test,

217 Bonferroni adjusted p-value 1*$10^{-40}$). By contrast, targets of marketed and clinical trial

218 drugs were significantly more conserved (the highest Bonferroni adjusted p-value was

219 1*$10^{-5}$ for oncology drug targets in phase 3). The trend held for targets for oncology

220 and non-oncology indications. $K_a/K_s$ ratios were inversely correlated with sequence

221 identity (Spearman rho -0.86) and similarity (-0.82) between human proteins and their

9

222    mouse orthologs. Conservation of protein sequence is considered a proxy for

223    conservation of biological function[26]. Also, 283 of 1,018 tissue-specific genes (27.8%)

224    compared to 2,719 of all 18,377 protein-coding genes (14.8%) did not have a unique

225    ortholog in mouse. Therefore, absence of a convenient animal model and gaps in

226    translation from animal research to clinical trials in humans may complicate

227    development of tissue-specific genes as drug targets.

228    We next examined selection pressure within the human species using a new metric,

229    that was recently developed by the ExAC consortium - probability of being loss-of-

230    function intolerant (pLI)[27]. Genes with high pLI have significantly lower observed than

231    expected frequencies of loss-of-function variants, indicating that deleterious variants

232    in these genes are selected out of the human population. Genes with pLI >= 0.9 are

233    considered loss-of-function intolerant and their "knockout" in humans implies "some

234    non-trivial survival or reproductive disadvantage"[27]. By contrast, genes with pLI <= 0.1

235    are considered loss-of-function tolerant[27]. In our analysis, tissue-specific genes were

236    enriched in loss-of-function tolerant and depleted of intolerant genes compared to all

237    protein-coding genes (**Fig. 4b**). By contrast, targets of oncology drugs were enriched

238    in loss-of-function intolerant (highest Bonferroni adjusted p-value $8*10^{-13}$ in phase 3)

239    and depleted of tolerant genes (highest Bonferroni adjusted p-value $3*10^{-9}$ for

240    marketed drugs). Targets of marketed non-oncology drugs had comparable

241    prevalence of loss-of-function tolerant and intolerant genes compared to all protein-

242    coding genes. The distributions of pLI confirmed that tissue-specific genes were more

243    likely to become targets for non-oncology drugs.

244    Genes with pLI >= 0.9 are more likely to be detected in genome-wide association

245    studies (GWAS)[27], and to attract attention as candidate drug targets through GWAS.

246    We investigated whether less conserved tissue-specific genes were less frequently

247   found in GWAS. GWAS variants are often located in intergenic regions and can be

248   mapped to candidate genes by proximity on the chromosome or through an

249   association between genotype of the GWAS variant and expression of a gene (eQTL).

250   Mapping through eQTL can highlight regulatory relationships in disease-relevant

251   tissues[28], and, consequently, is frequently used. The two types of mapping can

252   highlight different candidate genes[29,30], and require follow-up experiments to

253   determine causal genes. Tissue-specific genes were equally likely to be detected as

254   a nearest gene to a GWAS variant but 1.3 times less likely to be mapped from GWAS

255   to single-tissue cis-eQTLs than all protein-coding genes (Fisher test, Bonferroni

256   adjusted p-value $2*10^{-9}$). Interestingly, only mapping by proximity on chromosome

257   distinguished drug targets from all protein-coding genes (**Supplementary Fig. 2**). 111

258   of 392 (28.3%) of GWAS to eQTL relationships for tissue-specific genes were detected

259   in the corresponding tissues with highest expression. These results were not

260   surprising. Our definition of tissue-specificity allowed lower expression in other tissues.

261   Some tissues, including kidney cortex with 48 tissue-specific genes, had no eQTL data

262   due to insufficiently high number of samples and did not contribute to this analysis.

263   Also, approximately a third of GWAS to eQTL relationships can only be captured using

264   multiple tissues, while single-tissue analyses lack power to detect the associations[30].

265   Thus, tissue-specific genes were less likely to be highlighted as candidate targets if

266   investigators relied on the GWAS to eQTL approach.

267   In summary, underrepresentation of tissue-specific targets in early clinical trials could

268   be attributed to their primary relevance for non-oncology diseases and translational

269   challenges. Despite these challenges, tissue-specific genes were enriched among

270   targets of marketed drugs. We hypothesized that tissue-specificity was related to

271   efficacy and not only to safety.

11

## Tissue-specificity vs efficacy

Efficacy of a drug can be viewed as a combination of properties of the drug (e.g., potency, bioavailability, selectivity etc.) and properties of its intended target(s). Here, we focus on efficacy-related properties of the targets.

### Prevalence of disease genes

Drugs, that modulate targets with genetic evidence for a human disease, are less likely to fail in clinical trials for lack of efficacy[19]. Knowledge of human genetics can help to understand the biological function of the target, find target engagement biomarkers for clinical trials and estimate dose-response curves[31]. These factors can enhance the chances of a drug to succeed.

We compared the prevalence of OMIM[32] and Protein Truncating Variants *esc*aping nonsense mediated decay (PTVesc) genes[33] among tissue-specific genes and drug targets. OMIM genes have an entry in the Online Mendelian Inheritance in Man® Morbid Map data base[32], are well-known disease genes and are likely to be explored in target discovery. By constrast, PTVesc genes are an emerging class of candidate genes that can cause disease by gain-of-function mechanism. PTVesc genes are significantly depleted of genetic variants, that result in mRNA that escape nonsense-mediated decay and production of truncated proteins with altered function (e.g., *PNPLA3* and *APOL1*)[33]. Methods for detection of PTVesc are recently developed, so PTVesc genes are unlikely to be explored to the same extend as OMIM genes. Tissue-specific genes were enriched in both OMIM and PTVesc genes (**Fig. 5**). The outcomes of Fisher exact test for tissue-specific genes were 272 of 1,018 > 3,870 of 18,377, Bonferroni adjusted p-value $2*10^{-4}$ for OMIM genes and 159 of 1,018 > 1,913 of 18,377, Bonferroni adjusted p-value $4*10^{-6}$ for PTVesc genes. As expected, drug targets for oncology and non-oncology indications across all phases of clinical

297    development were enriched only in OMIM genes (**Fig. 5a**). By contrast, the prevalence

298    of PTVesc genes among drug targets did not significantly deviate from the overall

299    prevalence among protein-coding genes (**Fig. 5b**).

300    In total, 386 of 1,018 tissue-specific genes (37.9%) were OMIM genes or PTVesc

301    genes or both. Thus, tissue-specific genes were more likely to provide necessary

302    information for development of efficacious drugs through human genetics than protein-

303    coding genes overall.

304    **Network analysis**

305    The ability to spread perturbations through the cell and cause phenotypic changes is

306    a key property of drug targets, which is reflected by topological properties in protein-

307    protein interaction (PPI) networks[34]. We explain the network topology properties in

308    **Supplementary Fig. 3**. We recommend section 2 in[35] for a detailed explanation of the

309    relationship between network topology properties and spread of perturbations.

310    We performed network analysis on STRING v10.5[36] (**Supplementary Data 2**)

311    because tissue-specific proteins are well represented in this data base[37]. We included

312    three gene sets with known ability to affect phenotype as controls. Distribution of

313    network-topological properties of these genes should indicate if the PPI network

314    accurately reflects the ability of genes to spread perturbations and cause phenotype.

315    Essential genes cause cell death or hamper growth upon silencing in two human

316    cancer cell lines[38]. These genes serve as positive control for severe phenotypes.

317    OMIM genes cause disease and serve as positive control for less severe phenotypes.

318    Genes with rare homozygous loss-of-function rhLOF variants in three human cohorts

319    serve as negative no-phenotype controls (British-Pakistani, ExAC and Icelandic

320    individuals, Suppl. Table 2 from[39]). The human subjects come from the general

321    population and are assumed to be healthy, so loss of function of rhLOF genes is

13

322    assumed to be compensated. No association between presence of rhLOF genes and

323    rate of drug prescriptions and medical consultations has been confirmed in the British-

324    Pakistani cohort[39].

325    First, we investigated the sources of supporting evidence for PPIs (**Supplementary**

326    **Fig. 4**). Tissue-specific genes did not markedly differ from all protein-coding genes in

327    this respect. Each PPI had a score reflecting the amount of cumulative evidence

328    supporting existence of the interaction. Interestingly, non-oncology drug targets from

329    phase 1 to the market tended to have more high confidence PPIs than other gene

330    categories (**a**). PPIs for oncology drug targets and essential genes had more support

331    from co-expression across multiple experiments and tissues (**f**) and the experimental

332    evidence channel (**g**). PPIs of non-oncology drug targets tended to have more support

333    from pathway data bases (**h**). We concluded that indirect (*functional*) interactions were

334    important for non-oncology targets and kept both physical and functional interactions

335    for analysis. Most PPIs were supported by published scientific literature (**i**). The

336    number of reported PPIs and the number of published articles per gene were

337    correlated (Kendall tau b = 0.31), indicating a source of bias for network topology

338    analysis. The neighbourhood (**c**), fusion (**d**) and co-occurrence (**e**) channels provided

339    support for relatively few PPIs, consistent with primary relevance of these three

340    evidence channels for PPIs in Archaea and Bacteria[36].

341    We observed that most PPIs had low confidence scores even for drug targets (**b**).

342    Mora and Donaldson demonstrated that removing low confidence interaction does not

343    substantially improve the ability to discriminate drug targets based on their topological

344    properties[40]. Hence, we analysed the complete interaction set, but directly

345    incorporated the confidence in PPIs into the calculations and computed weighted

346    topological properties (see **Methods/Network analysis** for details). The calculations

14

347 were performed on the largest connected component including 19,574 proteins and

348 5,676,527 PPIs. The network diameter (unweighted) was 6. Topological properties of

349 the nodes accurately reflected their ability to spread perturbations through the network

350 (**Fig. 6** and **Table 1**). Distributions of centrality scores among rhLOF genes did not

351 significantly differ from the overall distributions among protein-coding genes (except

352 for slightly lower closeness centrality scores). Drug targets, OMIM genes and essential

353 genes had elevated centrality scores. Betweenness centrality was the only topological

354 property that could distinguish tissue-specific genes from all protein-coding genes

355 (**Table 1**). The trend was nominally significant but did not pass the correction for

356 multiple testing. Our results were consistent with the previous study on regulatory

357 networks, in which the Sonawane et al. applied a less stringent definition of tissue-

358 specificity and found that tissue-specific genes serve as "bottlenecks" on signaling

359 paths[41].

360 We further investigated which tissue-specific genes had high betweenness centrality

361 scores. The ten highest betweenness centrality scores were for genes encoding

362 hormones (insulin *INS*; glucagon *GCG*; *POMC* giving rise to adrenocorticotrophin and

363 lipotropin beta in the pituitary), other secreted proteins (albumin *ALB*; neuropeptide S

364 *NPS*; plasminogen *PLG*; *APOA1*, a major constituent of high density lipoprotein

365 cholesterol), rate limiting enzyme in synthesis of bile acids *CYP7A1*, mitochondrial

366 enzyme *FDXR* and electron transporter *FDX1* acting together in synthesis of steroid

367 hormones in the adrenal glands. Enzymes and secreted proteins, that were expressed

368 in a tissue-specific manner, had higher betweenness centrality scores than other

369 tissue-specific genes (**Supplementary Table 1**). These genes may have important

370 *functional* interactions and their modulation may cause effects outside of their tissue-

371 of-origin. For example, aliskiren fumarate inhibits the kidney-specific enzyme renin

15

372  *REN* that is part of the renin–angiotensin–aldosterone system, lowers blood pressure

373  and mitigates manifestations of hypertension in the whole body.

## Historic precedents to guide future applications

375  In total, only 100 of 1,018 (9.8%) tissue-specific genes were explored as targets of

376  marketed or clinical trial drugs. 284 of the remaining 918 (30.9%) tissue-specific genes

377  were classified as Tdark in the TCRD data base[42], i.e., were poorly researched with

378  unknown biological function. 529 of 918 (57.6%) showed some indication of

379  druggability by small molecule or antibody approaches (**Supplementary Fig. 5**). The

380  definition of druggability constantly expands, and targets that cannot be modulated

381  with small molecules or antibodies may be targeted by antisense oligonucleotides or

382  other approaches. Hence, the opportunities to identify novel drug targets among

383  tissue-specific genes were not exhausted. In the next subsections, we review how

384  tissue-specific genes have been used historically and highlight promising future

385  applications.

### Tissue-specific drug targets

387  Tissue-specific genes are targeted by drugs approved for diverse disease indications

388  (**Table 2**). Tissue-specific genes can be targeted by small molecules (e.g., *ACE* -

389  captopril), analogues of endogenous substances (*AVPR1B* – desmopressin acetate,

390  an analogue of vasopressin), antibodies (*TNF* - etanercept) and new modalities. We

391  and other researchers[9,17] demonstrated that tissue-specific genes are enriched in

392  OMIM genes. Hence, defective forms of tissue-specific genes causing rare monogenic

393  diseases are potential targets for genome editing with the emerging CRISPR/Cas9

394  technology (e.g., surfactant genes in surfactant deficiencies[43], *SERPINA1* in alpha-1

395  antitrypsin deficiency[44]).

16

396   Historically, tissue-specific genes were predominantly targets for non-oncology

397   diseases. However, tissue-specific genes also find applications in oncology (e.g.,

398   mitotane for endocrine therapy in inoperable adrenocortical carcinoma[45]) and as

399   targets for pharmacoenhancers. These application scenarios may be expanded in the

400   future.

401   The pharmacoenhancer Cobicistat is administered together with antiretroviral drugs

402   and inhibits cytochromes of *CYP3A* subfamily that degrade antiretroviral drugs

403   primarily in the liver and intestines. Cobicistat helps to maintain therapeutic

404   concentration of antiretroviral agents for a longer time whereby improving adherence

405   to therapy in HIV patients[46]. Drug-metabolizing enzymes such as the liver-specific

406   cytochromes *CYP1A2*, *CYP2D6* and *CYP2C9*[47] may be candidate targets for other

407   pharmcoenhancers to improve bioavailability or prolong action of the main drug.

408   Tissue-specific genes represent potential targets for antibody-based therapies (e.g.,

409   mammary gland specific transcription factor *ANKRD30A* for breast cancer[48]).

410   Promoters of tissue-specific genes can be used in oncolytic viral therapies to achieve

411   specific expression of the virus in the target tissue (e.g., urothelium-specific

412   adenovirus CG8840 with uroplakin 2 *UPK2* promoter for bladder cancer[49] and

413   prostate-specific antigen *KLK3* targeted adenovirus CG7870 for prostate cancer[50]). N-

414   acetylgalactosamine (GalNAc)-conjugated antisense oligonucleotide drugs bind to the

415   liver-specific *ASGR1* and enable targeted delivery to hepatocytes[51]. Similarly, tissue-

416   specific genes can be explored as targets for other targeted delivery approaches, and

417   not only in cancer.

418   Finally, *TNF* exemplifies a category of genes that dramatically change their expression

419   in disease and become specific to inflamed or cancerous tissue (represented by EBV-

420   transformed lymphoblastoid cell line and transformed fibroblasts in GTEx). Anti-TNF

421    drugs like etanercept are used to treat rheumatoid arthritis, psoriatic arthritis and

422    ankylosing spondylitis. *TNF* is 3.55 $\log_2$ fold (11.7 times on linear scale) higher in

423    synovial membranes of recently diagnosed patients with psoriathic arthritis, who are

424    naïve to anti-TNF treatment, than in healthy donors[52]. This example highlights the

425    importance of extensive tissue panels including both healthy and disease tissues such

426    as E-MTAB-3732[53] for target discovery in cancer and inflammatory diseases.

## Other applications

428    Tissue-specifically produced secreted proteins (or the corresponding recombinant

429    peptides, proteins and other synthetic analogs) are used as replacement therapy. The

430    best-known examples include hormone replacement therapies (insulin in type 1

431    diabetes, thyroid hormone in hypofunction of thyroid gland, oxytocin to induce labour,

432    etc.) and medication containing pancreatic or gastric enzymes to aid digestion (e.g.,

433    Creon). Other replacement therapies are in development. For example, FDA recently

434    approved the coagulation factor-albumin fusion protein Idelvion for congenital

435    complement factor IX *F9* deficiency while other complement replacement therapies

436    are in clinical trials[54]. In our analysis, 268 of 918 tissue-specific genes, that were not

437    yet explored as targets of marketed or clinical trial drugs, encoded secreted proteins,

438    which may represent opportunities for new replacement therapies (e.g., hormones and

439    their combinations for treatment of type 2 diabetes and obesity[55], artificial saliva with

440    recombinant lysozyme in treatment of xerostomia and Sjögren's syndrome[56]).

441    Tissue-specific genes with protein products entering the bloodstream can serve as

442    biomarkers to monitor the state or function of a tissue. For example, *KLK3*, better

443    known as prostate cancer antigen, has prostate-specific expression, enters

444    bloodstream, displays elevated levels in prostate cancer and benign hyperplasia of the

445    prostate and is used as a pre-screening test for prostate cancer[57]. Tissue-specific

18

446 biomarkers indicating tissue damage have several conceptual advantages over

447 conventional laboratory tests including higher specificity[58]. For example, circulating

448 proteins with liver-specific expression are evaluated as markers of acetaminophen

449 induced hepatotoxicity[59]. Biomarkers discussed in literature tended to conform to more

450 liberal definitions of tissue-specificity, which may be sufficiently stringent for successful

451 development of biomarkers (e.g., *RBP4*, evaluated in[59], satisfied constraint $x = 5$ in

452 liver and had lower expression in adipose tissues and pituitary, *NPPB* encoding NT-

453 pro-BNP satisfied $x = 4$ in heart atrial appendage with high expression in some left

454 ventricle samples).

## Non-coding genes

456 We focused on protein-coding genes, but long non-coding RNA (lncRNA) also tend to

457 be expressed in tissue- and cell-type specific manner[60]. We identified 2,113 long non-

458 coding RNAs (**Supplementary Data 3**), from which 77 were tissue-specific at $x = 6$

459 and confirmed polyadenylated, i.e., reliably detected with the GTEx RNA sequencing

460 protocol. Long non-coding RNAs have potential as drug targets[61]. For example,

461 prostate-specific lncRNA *PCGEM1* is candidate target in prostate cancer[62].

# DISCUSSION

463 We conducted a retrospective analysis of tissue-specific genes compared to drug

464 targets in all phases of clinical development. Targets of phase 1 drugs reflect the most

465 recent research. By contrast, targets of marketed drugs have undergone at least a

466 decade in preclinical and clinical development and reflect older research.

467 Theoretically, overrepresentation of tissue-specific targets on the market and

468 depletion in phase 1 could reflect a historic shift in target selection paradigms.

469 However, Rouillard and colleagues[20] studied phase 3 drugs (projects with comparable

470 "age") and demonstrated that drugs modulating tissue-specific targets are more likely

471    to succeed in phase 3 and gain regulatory approval. Thus, the data presented in Fig.

472    2 and Supplementary Fig. 1 do not merely represent a historic trend. We are justified

473    to state that drugs modulating tissue-specific targets are indeed more likely to progress

474    in the clinic. We observed that tissue-specific genes, satisfying x = 6, were 1.8 times

475    more likely to become targets of marketed drugs and 2.3 times more likely to become

476    targets of marketed non-oncology drugs than protein-coding genes overall. Our

477    findings were consistent with the previous studies[8,9].

478    Success of tissue-specific genes as drug targets may be due to a complex

479    combination of factors. Good understanding of target biology is essential for target-

480    based drug discovery. Tissue-specific genes are involved in specialized tissue

481    functions and human diseases, in which genetic evidence can provide the necessary

482    supporting information for development of new drugs. By contrast, broadly expressed

483    genes may have diverging functions and protein isoforms with distinct subcellular

484    localization in different tissues[7]. The biology of tissue-specific genes may be less

485    complex to study than the biology of broadly expressed genes, which may contribute

486    to the development of efficacious drugs. The ability to spread perturbations within and

487    outside the tissue-of origin may be a direct contributing factor to efficacy of drugs

488    modulating tissue-specific genes with high betweenness centrality scores, especially

489    enzymes and genes encoding secreted proteins. Narrow expression profile decreases

490    the probability of side effects in non-intended target tissues. Tissue-specific genes are

491    depleted in loss-of-function intolerant genes, have average number of neighbors and

492    are located in interactome regions with average connectivity as indicated by

493    distributions of strength and eigenvector centrality scores. These properties indicate

494    improved safety, as they are in sharp contrast to oncology targets, whose modulation

495    can cause severe side effect.

496    Feasibility is another crucial component to the success of a drug program. Murine

497    models are commonly used to assess toxicity and for early efficacy studies *in vivo* and

498    can serve as a 'filter' to make stop/go decisions for a drug project. The targets of drugs

499    from phase 1 to the market are biased towards evolutionary conserved genes. We

500    suggest using humanized mouse models or other non-rodent species with smaller

501    evolutionary distance to humans to overcome translational challenges and enable the

502    development of tissue-specific genes and other less conserved biological entities like

503    long non-coding RNAs[63] as drug targets.

# MATERIALS AND METHODS

## Gene expression

506    Gene-level RPKM values were downloaded from The Genotype-Tissue Expression

507    Consortium[30] (https://gtexportal.org/home/, release 6). The per-tissue mean RPKM for

508    each gene was subjected to Z-transformation across tissues and then to a second Z-

509    transformation across genes to bring all Z-scores to the same scale. We identified

510    18,377 protein-coding genes and 2,113 long non-coding RNA with HGNC approved

511    gene symbol[64]. The non-alternative loci data set was obtained from the HGNC

512    Database (www.genenames.org, 30.08.2017).

## Drug targets

514    Mechanism-of-action targets of marketed and clinical trial drugs, disease indications

515    and year of first approval for marketed drugs were extracted from ChEMBL version

516    23[65]. Drugs were classified as phase 1, 2, 3 or marketed drugs based on the maximal

517    phase they reached in clinical trials. Disease indications were mapped to Disease

518    Ontology[66]. Proteins were classified as oncology or non-oncology targets based on

519 parent terms in Disease Ontology. If a protein was targeted by at least 1 oncology

520 drug, it was considered an oncology target.

521 ## Meta-data

522 Example compounds with exact $K_i$ or $IC_{50}$ activity values against human proteins,

523 measured in assays with direct interaction and the highest confidence score=9, were

524 retrieved from ChEMBL v23[65]. Mapping from ENSEMBL identifiers to PDB and

525 polyadenylation data were obtained from GENECODE consortium[67] version 27.

526 Mapping to enzyme EC numbers, Uniprot and NCBI Gene (Entrez) identifiers were

527 extracted from the HGNC non-alternative loci data set[64]. Target Development Level

528 (TDL) was retrieved from TCRD version 4.6.2[42]. Subcellular localization and protein

529 family information were obtained from UniProt/SwissProt[68]. Probabilities of being loss-

530 of-function intolerant (pLI) were retrieved from Supplementary Data of the ExAC

531 consortium flagship publication[27]. Associations with Mendelian diseases were

532 retrieved from OMIM Morbid Map[32] (copyright John Hopkins University, AstraZeneca

533 purchased license JHU agreement number A30699 and reference number C03746).

534 We included only binary indicator variables (has/does not have an entry in the Morbid

535 Map). Number of PubMed-indexed articles linked to each gene was retrieved from

536 NCBI Gene[69] https://www.ncbi.nlm.nih.gov/gene/ on the 02.01.2018. Human to mouse

537 orthologs, $K_a/K_s$ ratios and percentages of sequence identity and similarity were

538 extracted from ENSEMBL Compara[70] version 91. The lists of essential genes[38],

539 PTVesc[33] and rhLOF[39] genes were obtained from supplementary data of the

540 respective publications. Biological function of the genes was described according to

541 the NCBI Gene/RefSeq summary[71] unless explicitly indicated otherwise.

542 ## Network analysis

543  Human protein-protein interaction network was downloaded from STRING v 10.5 (file

544  9606.protein.links.detailed.v10.5.tsv)[36]. Topological properties were calculated with

545  igraph[72] version 1.2.1. Weighted k-shell decomposition was computed as described

546  in[73]. Combined evidence scores were used as edge weights for strength, eigenvector

547  centrality and k-shell calculations, i.e., the overall 'influence' of a node was

548  proportional to the number of its neighbors combined with confidence in its PPIs. Edge

549  weights were taken as (1 – combined evidence score) for centrality measures based

550  on shortest paths, i.e., shortest paths were the 'least uncertainty paths'.

## Mapping from GWAS to candidate genes

552  Genetic associations were obtained from GWAS Catalog[74] (data set: all

553  associations v1.0, access: 30.08.2017, https://www.ebi.ac.uk/gwas/). Coordinates of

554  genetic variants (SNPs) were mapped from genome assembly GRCh38 to GRCh37

555  by SNP identifiers (rsids) in 1000 genomes phase 3[75]. Proxy SNPs with r2 >= 0.8 were

556  identified within 50 kilobasepairs in the CEU population using *--hap-r2*-positions

557  command with vcftools[76] version 0.1.13. GWAS variants and their proxy SNPs were

558  mapped to significant single-tissue eQTL from GTEx[30] version 7.

## Statistical analysis

560  We applied Fisher exact test for count data because sample sizes were small in some

561  instances (e.g., 3 tissue-specific targets in phase 1 at x = 6) and to be consistent in

562  other analyses. Mann-Whitney U test was used to test differences between groups for

563  continuous variables. Wilcoxon test with explicit handling of tied values in

564  exactRankTests[77] version 0.8-29 was used to test differences in year of first approval

565  in Fig. 3c. Tests for enrichment or depletion were one-tailed, other tests were two-

566  tailed. Bonferroni correction for multiple testing was applied as appropriate. P-values

23

567 < 0.05 were considered significant. Statistical analyses were summarized in

568 **Supplementary Information 1**. Figures were generated with ggplot2[78] version 3.0.0,

569 viridis[79] version 0.5.1, VennDiagram[80] version 1.6.20 and UpSetR[81] version 1.3.3.

570 Analyses were performed in R[82] version 3.4.1.

571 ## Data availability statement

572 All data, that were generated in this study, are provided as Supplementary data sets.

573 Annotated Z-score tables for protein-coding genes including the tissue-specific gene

574 and drug target subsets are provided in **Supplementary Data 1**. Network topology

575 properties are provided in **Supplementary Data 2**. Long non-coding RNAs are listed

576 in **Supplementary Data 3**. Columns, that were used as input data for figures, are

577 labelled within each supplementary data set. Summary-level data (counts and

578 percentages) behind figures are included in the **Supplementary Information 1.**

579 Source data for Fig. 4a can be retrieved directly from Ensembl Compara[70] v 91.

580

581 ## REFERENCES

582 1    Debouck, C. & Goodfellow, P. N. DNA microarrays in drug discovery and

583      development. *Nat Genet* **21**, 48-50, doi:10.1038/4475 (1999).

584 2    Gashaw, I., Ellinghaus, P., Sommer, A. & Asadullah, K. What makes a good

585      drug    target?    *Drug    Discov    Today*    **16**,    1037-1043,

586      doi:10.1016/j.drudis.2011.09.007 (2011).

587 3    Xu, H. *et al.* Learning the drug target-likeness of a protein. *Proteomics* **7**, 4255-

588      4263, doi:10.1002/pmic.200700062 (2007).

589  4  Yao, L. & Rzhetsky, A. Quantitative systems-level determinants of human genes targeted by successful drugs. *Genome Res* **18**, 206-213, doi:10.1101/gr.6888208 (2008).

592  5  Kim, B., Jo, J., Han, J., Park, C. & Lee, H. In silico re-identification of properties of drug target proteins. *BMC Bioinformatics* **18**, 248, doi:10.1186/s12859-017-1639-3 (2017).

595  6  Emig, D. & Albrecht, M. Tissue-specific proteins and functional implications. *J Proteome Res* **10**, 1893-1903, doi:10.1021/pr101132h (2011).

597  7  Uhlen, M. *et al.* Proteomics. Tissue-based map of the human proteome. *Science* **347**, 1260419, doi:10.1126/science.1260419 (2015).

599  8  Dezso, Z. *et al.* A comprehensive functional analysis of tissue specificity of human gene expression. *BMC Biol* **6**, 49, doi:10.1186/1741-7007-6-49 (2008).

601  9  Yang, L. *et al.* Comparative analysis of housekeeping and tissue-selective genes in human based on network topologies and biological properties. *Mol Genet Genomics* **291**, 1227-1241, doi:10.1007/s00438-016-1178-z (2016).

604  10  Kumar, V., Sanseau, P., Simola, D. F., Hurle, M. R. & Agarwal, P. Systematic Analysis of Drug Targets Confirms Expression in Disease-Relevant Tissues. *Sci Rep* **6**, 36205, doi:10.1038/srep36205 (2016).

607  11  Liu, X., Yu, X., Zack, D. J., Zhu, H. & Qian, J. TiGER: a database for tissue-specific gene expression and regulation. *BMC Bioinformatics* **9**, 271, doi:10.1186/1471-2105-9-271 (2008).

610  12  Xiao, S. J., Zhang, C., Zou, Q. & Ji, Z. L. TiSGeD: a database for tissue-specific genes. *Bioinformatics* **26**, 1273-1275, doi:10.1093/bioinformatics/btq109 (2010).

613    13    Yang, X. *et al.* VeryGene: linking tissue-specific genes to diseases, drugs, and

614          beyond for knowledge discovery. *Physiol Genomics* **43**, 457-460,

615          doi:10.1152/physiolgenomics.00178.2010 (2011).

616    14    Kim, P. *et al.* TissGDB: tissue-specific gene database in cancer. *Nucleic Acids*

617          *Res* **46**, D1031-D1038, doi:10.1093/nar/gkx850 (2018).

618    15    DiMasi, J. A., Grabowski, H. G. & Hansen, R. W. Innovation in the

619          pharmaceutical industry: New estimates of R&D costs. *J Health Econ* **47**, 20-

620          33, doi:10.1016/j.jhealeco.2016.01.012 (2016).

621    16    Oprea, T. I. *et al.* Unexplored therapeutic opportunities in the human genome.

622          *Nat Rev Drug Discov* **17**, 317-332, doi:10.1038/nrd.2018.14 (2018).

623    17    Winter, E. E., Goodstadt, L. & Ponting, C. P. Elevated rates of protein secretion,

624          evolution, and disease among tissue-specific genes. *Genome Res* **14**, 54-61,

625          doi:10.1101/gr.1924004 (2004).

626    18    Nelson, M. R. *et al.* The support of human genetic evidence for approved drug

627          indications. *Nat Genet* **47**, 856-860, doi:10.1038/ng.3314 (2015).

628    19    Cook, D. *et al.* Lessons learned from the fate of AstraZeneca's drug pipeline: a

629          five-dimensional framework. *Nat Rev Drug Discov* **13**, 419-431,

630          doi:10.1038/nrd4309 (2014).

631    20    Rouillard, A. D., Hurle, M. R. & Agarwal, P. Systematic interrogation of diverse

632          Omic data reveals interpretable, robust, and generalizable transcriptomic

633          features of clinically successful therapeutic targets. *PLoS Comput Biol* **14**,

634          e1006142, doi:10.1371/journal.pcbi.1006142 (2018).

635    21    Liu, Y., Beyer, A. & Aebersold, R. On the Dependency of Cellular Protein Levels

636          on mRNA Abundance. *Cell* **165**, 535-550, doi:10.1016/j.cell.2016.03.014

637          (2016).

26

638    22    Edfors, F. *et al.* Gene-specific correlation of RNA and protein levels in human
639          cells and tissues. *Mol Syst Biol* **12**, 883, doi:10.15252/msb.20167144 (2016).
640    23    Schafer, S. *et al.* Translational regulation shapes the molecular landscape of
641          complex disease phenotypes. *Nat Commun* **6**, 7200, doi:10.1038/ncomms8200
642          (2015).
643    24    Lv, W. *et al.* The drug target genes show higher evolutionary conservation than
644          non-target genes. *Oncotarget* **7**, 4961-4971, doi:10.18632/oncotarget.6755
645          (2016).
646    25    Hurst, L. D. The Ka/Ks ratio: diagnosing the form of sequence evolution. *Trends*
647          *Genet* **18**, 486 (2002).
648    26    Joshi, T. & Xu, D. Quantitative assessment of relationship between sequence
649          similarity and function similarity. *BMC Genomics* **8**, 222, doi:10.1186/1471-
650          2164-8-222 (2007).
651    27    Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans.
652          *Nature* **536**, 285-291, doi:10.1038/nature19057 (2016).
653    28    Hauberg, M. E. *et al.* Large-Scale Identification of Common Trait and Disease
654          Variants Affecting Gene Expression. *Am J Hum Genet* **101**, 157,
655          doi:10.1016/j.ajhg.2017.06.003 (2017).
656    29    Nica, A. C. *et al.* Candidate causal regulatory effects by integration of
657          expression QTLs with complex trait genetic associations. *PLoS Genet* **6**,
658          e1000895, doi:10.1371/journal.pgen.1000895 (2010).
659    30    GTEx Consortium. Human genomics. The Genotype-Tissue Expression
660          (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648-
661          660, doi:10.1126/science.1262110 (2015).

662    31    Plenge, R. M., Scolnick, E. M. & Altshuler, D. Validating therapeutic targets
663          through human genetics. *Nat Rev Drug Discov* **12**, 581-594,
664          doi:10.1038/nrd4051 (2013).

665    32    McKusick-Nathans Institute of Genetic Medicine & Johns Hopkins University
666          (Baltimore, MD). Online Mendelian Inheritance in Man, OMIM®. World Wide
667          Web URL: https://omim.org/.  (2017).

668    33    Coban-Akdemir, Z. *et al.* Identifying Genes Whose Mutant Transcripts Cause
669          Dominant Disease Traits by Potential Gain-of-Function Alleles. *Am J Hum
670          Genet* **103**, 171-187, doi:10.1016/j.ajhg.2018.06.009 (2018).

671    34    Perez-Lopez, A. R. *et al.* Targets of drugs are generally, and targets of drugs
672          having side effects are specifically good spreaders of human interactome
673          perturbations. *Sci Rep* **5**, 10182, doi:10.1038/srep10182 (2015).

674    35    Pei, S., Morone, F. & Makse, H. A. in *Complex Spreading Phenomena in Social
675          Systems*    125-148 (Springer, 2018).

676    36    Szklarczyk, D. *et al.* The STRING database in 2017: quality-controlled protein-
677          protein association networks, made broadly accessible. *Nucleic Acids Res* **45**,
678          D362-D368, doi:10.1093/nar/gkw937 (2017).

679    37    Lopes, T. J. *et al.* Tissue-specific subnetworks and characteristics of publicly
680          available human protein interaction databases. *Bioinformatics* **27**, 2414-2421,
681          doi:10.1093/bioinformatics/btr414 (2011).

682    38    Blomen, V. A. *et al.* Gene essentiality and synthetic lethality in haploid human
683          cells. *Science* **350**, 1092-1096, doi:10.1126/science.aac7557 (2015).

684    39    Narasimhan, V. M. *et al.* Health and population effects of rare gene knockouts
685          in adult humans with related parents. *Science* **352**, 474-477,
686          doi:10.1126/science.aac8624 (2016).

28

687  40  Mora, A. & Donaldson, I. M. Effects of protein interaction data integration,

688      representation and reliability on the use of network properties for drug target

689      prediction. *BMC Bioinformatics* **13**, 294, doi:10.1186/1471-2105-13-294 (2012).

690  41  Sonawane, A. R. *et al.* Understanding Tissue-Specific Gene Regulation. *Cell*

691      *Rep* **21**, 1077-1088, doi:10.1016/j.celrep.2017.10.001 (2017).

692  42  Nguyen, D. T. *et al.* Pharos: Collating protein information to shed light on the

693      druggable    genome.    *Nucleic    Acids    Res*    **45**,    D995-D1002,

694      doi:10.1093/nar/gkw1072 (2017).

695  43  Alapati, D. & Morrisey, E. E. Gene Editing and Genetic Lung Disease. Basic

696      Research Meets Therapeutic Application. *Am J Respir Cell Mol Biol* **56**, 283-

697      290, doi:10.1165/rcmb.2016-0301PS (2017).

698  44  Loring, H. S. & Flotte, T. R. Current status of gene therapy for alpha-1

699      antitrypsin    deficiency.    *Expert    Opin    Biol    Ther*    **15**,    329-336,

700      doi:10.1517/14712598.2015.978854 (2015).

701  45  Veytsman, I., Nieman, L. & Fojo, T. Management of endocrine manifestations

702      and the use of mitotane as a chemotherapeutic agent for adrenocortical

703      carcinoma. *J Clin Oncol* **27**, 4619-4629, doi:10.1200/JCO.2008.17.2775

704      (2009).

705  46  Xu, L. *et al.* Cobicistat (GS-9350): A Potent and Selective Inhibitor of Human

706      CYP3A as a Novel Pharmacoenhancer. *ACS Med Chem Lett* **1**, 209-213,

707      doi:10.1021/ml1000257 (2010).

708  47  Wilkinson, G. R. Drug metabolism and variability among patients in drug

709      response. *N Engl J Med* **352**, 2211-2221, doi:10.1056/NEJMra032424 (2005).

710  48  Seil, I. *et al.* The differentiation antigen NY-BR-1 is a potential target for
711       antibody-based therapies in breast cancer. *Int J Cancer* **120**, 2635-2642,
712       doi:10.1002/ijc.22620 (2007).

713  49  Zhang, J. *et al.* Identification of human uroplakin II promoter and its use in the
714       construction of CG8840, a urothelium-specific adenovirus variant that
715       eliminates established bladder tumors in combination with docetaxel. *Cancer*
716       *Res* **62**, 3743-3750 (2002).

717  50  Small, E. J. *et al.* A phase I trial of intravenous CG7870, a replication-selective,
718       prostate-specific antigen-targeted oncolytic adenovirus, for the treatment of
719       hormone-refractory, metastatic prostate cancer. *Mol Ther* **14**, 107-117,
720       doi:10.1016/j.ymthe.2006.02.011 (2006).

721  51  Tanowitz, M. *et al.* Asialoglycoprotein receptor 1 mediates productive uptake of
722       N-acetylgalactosamine-conjugated and unconjugated phosphorothioate
723       antisense oligonucleotides into liver hepatocytes. *Nucleic Acids Res* **45**, 12388-
724       12400, doi:10.1093/nar/gkx960 (2017).

725  52  Dolcino, M. *et al.* Gene Expression Profiling in Peripheral Blood Cells and
726       Synovial Membranes of Patients with Psoriatic Arthritis. *PLoS One* **10**,
727       e0128262, doi:10.1371/journal.pone.0128262 (2015).

728  53  Torrente, A. *et al.* Identification of Cancer Related Genes Using a
729       Comprehensive Map of Human Gene Expression. *PLoS One* **11**, e0157484,
730       doi:10.1371/journal.pone.0157484 (2016).

731  54  Escobar, M. A. Advances in the treatment of inherited coagulation disorders.
732       *Haemophilia* **19**, 648-659, doi:10.1111/hae.12137 (2013).

733   55   Sadry, S. A. & Drucker, D. J. Emerging combinatorial hormone therapies for the

734          treatment of obesity and T2DM. *Nat Rev Endocrinol* **9**, 425-433,

735          doi:10.1038/nrendo.2013.47 (2013).

736   56   Tenovuo, J. Clinical applications of antimicrobial host proteins lactoperoxidase,

737          lysozyme and lactoferrin in xerostomia: efficacy and safety. *Oral Dis* **8**, 23-29

738          (2002).

739   57   De Angelis, G., Rittenhouse, H. G., Mikolajczyk, S. D., Blair Shamel, L. &

740          Semjonow, A. Twenty Years of PSA: From Prostate Antigen to Tumor Marker.

741          *Rev Urol* **9**, 113-123 (2007).

742   58   Muller, P. Y. & Dieterle, F. Tissue-specific, non-invasive toxicity biomarkers:

743          translation from preclinical safety assessment to clinical safety monitoring.

744          *Expert Opin Drug Metab Toxicol* **5**, 1023-1038,

745          doi:10.1517/17425250903114174 (2009).

746   59   Hu, Z. *et al.* Quantitative liver-specific protein fingerprint in blood: a signature

747          for hepatotoxicity. *Theranostics* **4**, 215-228, doi:10.7150/thno.7868 (2014).

748   60   Hon, C. C. *et al.* An atlas of human long non-coding RNAs with accurate 5'

749          ends. *Nature* **543**, 199-204, doi:10.1038/nature21374 (2017).

750   61   Matsui, M. & Corey, D. R. Non-coding RNAs as drug targets. *Nat Rev Drug*

751          *Discov* **16**, 167-179, doi:10.1038/nrd.2016.117 (2017).

752   62   Hung, C. L. *et al.* A long noncoding RNA connects c-Myc to tumor metabolism.

753          *Proc Natl Acad Sci U S A* **111**, 18697-18702, doi:10.1073/pnas.1415669112

754          (2014).

755   63   Hezroni, H. *et al.* Principles of long noncoding RNA evolution derived from

756          direct comparison of transcriptomes in 17 species. *Cell Rep* **11**, 1110-1122,

757          doi:10.1016/j.celrep.2015.04.023 (2015).

758    64    Yates, B. *et al.* Genenames.org: the HGNC and VGNC resources in 2017.

759          *Nucleic Acids Res* **45**, D619-D625, doi:10.1093/nar/gkw1033 (2017).

760    65    Gaulton, A. *et al.* The ChEMBL database in 2017. *Nucleic Acids Res* **45**, D945-

761          D954, doi:10.1093/nar/gkw1074 (2017).

762    66    Kibbe, W. A. *et al.* Disease Ontology 2015 update: an expanded and updated

763          database of human diseases for linking biomedical knowledge through disease

764          data. *Nucleic Acids Res* **43**, D1071-1078, doi:10.1093/nar/gku1011 (2015).

765    67    Harrow, J. *et al.* GENCODE: the reference human genome annotation for The

766          ENCODE Project. *Genome Res* **22**, 1760-1774, doi:10.1101/gr.135350.111

767          (2012).

768    68    The UniProt Consortium. UniProt: the universal protein knowledgebase.

769          *Nucleic Acids Res* **46**, 2699, doi:10.1093/nar/gky092 (2018).

770    69    NCBI Resource Coordinators. Database resources of the National Center for

771          Biotechnology Information. *Nucleic Acids Res* **46**, D8-D13,

772          doi:10.1093/nar/gkx1095 (2018).

773    70    Herrero, J. *et al.* Ensembl comparative genomics resources. *Database (Oxford)*

774          **2016**, doi:10.1093/database/bav096 (2016).

775    71    O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current

776          status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**,

777          D733-745, doi:10.1093/nar/gkv1189 (2016).

778    72    Csardi, G. & Nepusz, T. The igraph software package for complex network

779          research. *InterJournal, Complex Systems* **1695**, 1-9 (2006).

780    73    Garas, A., Schweitzer, F. & Havlin, S. A k-shell decomposition method for

781          weighted networks. *New Journal of Physics* **14**, 083030 (2012).

782    74    MacArthur, J. *et al.* The new NHGRI-EBI Catalog of published genome-wide

783          association studies (GWAS Catalog). *Nucleic Acids Res* **45**, D896-D901,

784          doi:10.1093/nar/gkw1133 (2017).

785    75    1000 Genomes Project Consortium *et al.* A global reference for human genetic

786          variation. *Nature* **526**, 68-74, doi:10.1038/nature15393 (2015).

787    76    Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**,

788          2156-2158, doi:10.1093/bioinformatics/btr330 (2011).

789    77    Hothorn, T. & Hornik, K. exactRankTests: Exact Distributions for Rank and

790          Permutation    Tests.    https://CRAN.R-project.org/package=exactRankTests

791          (2017).

792    78    Wickham, H. *ggplot2: elegant graphics for data analysis*.  (Springer, 2016).

793    79    Garnier, S. viridis: Default Color Maps from 'matplotlib'. https://CRAN.R-

794          project.org/package=viridis (2018).

795    80    Chen, H. VennDiagram: Generate High-Resolution Venn and Euler Plots.

796          https://CRAN.R-project.org/package=VennDiagram (2018).

797    81    Conway, J. R., Lex, A. & Gehlenborg, N. UpSetR: an R package for the

798          visualization of intersecting sets and their properties. *Bioinformatics* **33**, 2938-

799          2940, doi:10.1093/bioinformatics/btx364 (2017).

800    82    R Core Team. R: A language and environment for statistical computing. (R

801          Foundation for Statistical Computing, Vienna, Austria., 2018).

802

803

804

805

806

## Author contribution statement

808  Conceptualization: MR, MH. Formal data analysis: MR. Writing: MR, MH.

## Conflicting interest

810  MR is a contractor to AstraZeneca. MH is employed by AstraZeneca. AstraZeneca

811  provided support to the authors in form of salaries, but had no role in conceptualization

812  of the study, data collection, analysis, interpretation and writing.

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831 **Table 1. Betweenness centrality scores in STRING v 10.5.** Betweenness centrality

832 scores are displayed separately in tabular form due to skewed distributions. IQR

833 stands for interquartile range. P-values are from two tailed Mann-Whitney U test

834 between the gene categories and all protein-coding genes (marked as 'Reference').

| Group | Median (IQR) | Nominal p | Bonferroni p |
|---|---|---|---|
| All proteins, N=19,574 | 14 (0-19,548) | Reference | Reference |
| Tissue-specific (x=6), N=1,004 | 32.1 (0-19,536.3) | $7.3*10^{-3}$ | 0.088 |
| Essential, N=1,713 | 29,414.2 (40-130,570) | $8.9*10^{-179}$ | $1.1*10^{-177}$ |
| OMIM, N=3,844 | 14,754.6 (7-59,036.2) | $2.6*10^{-169}$ | $3.1*10^{-168}$ |
| rhLOF, N=107 | 8 (0-2,241.5) | 0.2 | 1 |
| Marketed, oncology, N=211 | 48,734.4 (11,077.1-201,844.9) | $2.0*10^{-46}$ | $2.4*10^{-45}$ |
| Marketed, non-oncology, N=477 | 1,331.4 (13-44,616.7) | $5.5*10^{-23}$ | $6.6*10^{-22}$ |
| Phase 3, oncology, N=146 | 69,496.3 (12,503.2-357,208.5) | $3.6*10^{-33}$ | $4.3*10^{-32}$ |
| Phase 3, non-oncology, N=266 | 8,921.6 (38.5-59,684.4) | $2.7*10^{-21}$ | $3.3*10^{-20}$ |
| Phase 2, oncology, N=239 | 58,807 (595.9-375,702) | $9.4*10^{-46}$ | $1.1*10^{-44}$ |
| Phase 2, non-oncology, N=315 | 1,996.8 (9-53,038) | $2.5*10^{-15}$ | $3.0*10^{-14}$ |
| Phase 1, oncology, N=253 | 51,848.7 (5,140.7-222,555.3) | $6.8*10^{-46}$ | $8.1*10^{-45}$ |
| Phase 1, non-oncology, N=79 | 19,390.9 (53.2-7,7431.6) | $1.4*10^{-7}$ | $1.7*10^{-6}$ |

835

836

837

838

839

840

841

842

843

**Table 2. Examples of tissue-specific (x = 6) targets of marketed drugs.**

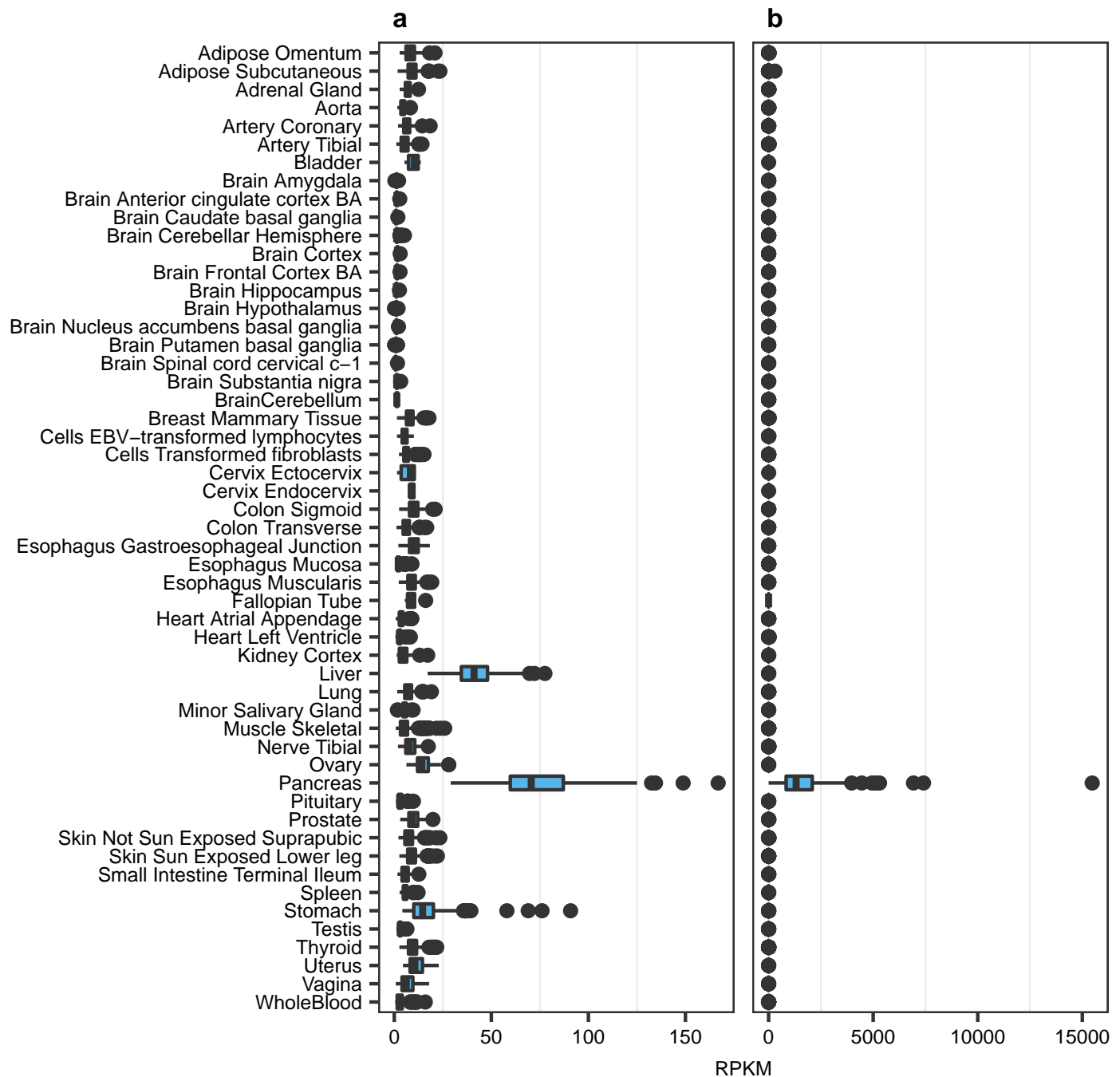| Target genes (tissue with highest expression) | Drug and indication |
|---|---|
| *CYP11A1*, *CYP11B1* (Adrenal gland) | Mitotane (adrenocortical carcinoma) |
| *TNF* (EBV-transformed lymphocytes) | Etanercept (autoimmune disease) |
| *SLC5A2* (Kidney) | Empagliflozin (type 2 diabetes) |
| *SLC22A6*, *SLC22A8* and *SLC22A11* (Kidney) | Probenecid (gout) |
| *CPS1* (Liver) | Carglumic acid (hyperammonaemia) |
| *CYP3A7* (Liver) | Cobicistat (HIV infection) |
| *CALCR* (Hypothalamus) | Calcitonin (hypercalcemia, osteoporosis) |
| *PNLIP* (Pancreas), *LIPF* (Stomach) | Orlistat (obesity) |
| *AVPR1B* (Pituitary) | Desmopressin acetate (diabetes Insipidus) |
| *CHRNA1*, *CHRND*, *CHRNG* (Skeletal muscle) | Atracurium besilate (myorelaxant in surgery) |
| *ACE* (Terminal ileum) | Captopril (hypertension) |
| *ATP4A*, *ATP4B* (Stomach) | Omeprazole (peptic ulcers) |
| *SLC6A3* (Substantia nigra) | Armodafinil (sleep disorders) |
| *TSHR* (Thyroid) | Thyrotropin alpha (thyroid cancer) |
| *CSF3R* (Whole blood) | Filgrastim (neutropenia) |

845

**Fig. 1. Examples of genes satisfying the most liberal (x = 2) and the most stringent (x = 10) definitions of tissue-specificity in pancreas.** **a** Amino acid transporter *SLC43A1* satisfies constraint x = 2 and does not satisfy more stringent definitions. **b** Hormone insulin *INS* satisfies the most stringent definition x = 10 and consequently satisfies all more liberal definitions x = 2 to 9. The box plots show normalized expression levels in RPKM across 53 human tissues in GTEx release 6 (https://gtexportal.org/home/).
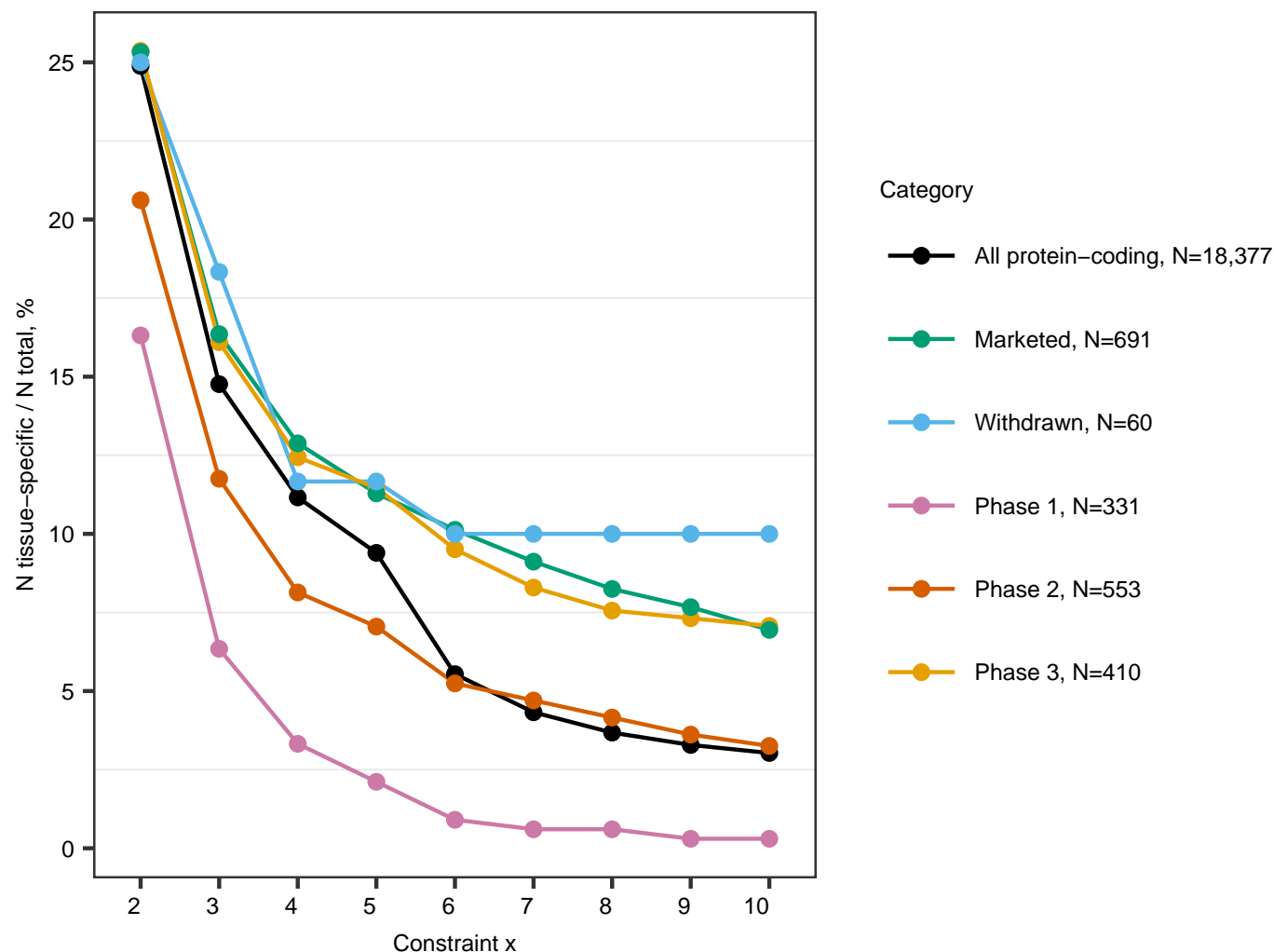
**Fig. 2. Prevalence of tissue-specific targets increased from phase 1 to the market.** Percentages of tissue-specific genes among targets of drugs in each phase of clinical development were plotted in comparison to the "background" distribution among all protein coding-genes (black line). Tissue-specificity was defined at nine increasingly stringent constraints x = 2 to 10 as illustrated in Fig. 1.
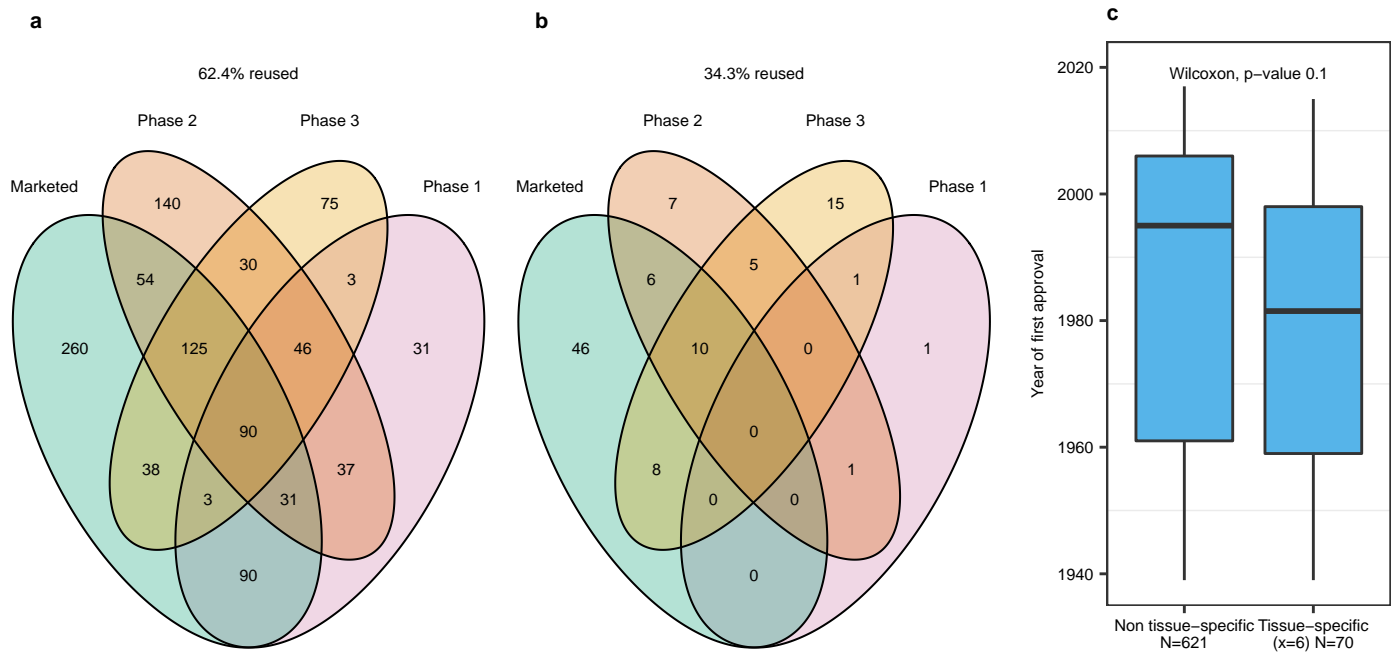
**Fig. 3. Tissue-specific targets, satisfying constraint x = 6, represented a less frequently reused (b vs a) and older (c) subset of targets of marketed drugs. a** Reuse of all targets of marketed drugs by other drugs in clinical trials. **b** Reuse of tissue-specific ($x = 6$) targets of marketed drugs by other drugs in clinical trials. Venn diagrams depict the number of targets in each phase of clinical development. Overlapping areas contain genes that are targeted by several drugs in different phases of development. **c** Year of regulatory approval by FDA or another agency of the first drugs modulating non-tissue-specific targets compared to tissue-specific ($x = 6$) targets of marketed drugs. For example, carglumic acid was the first marketed drug modulating *CPS1* and it was approved in 2010.
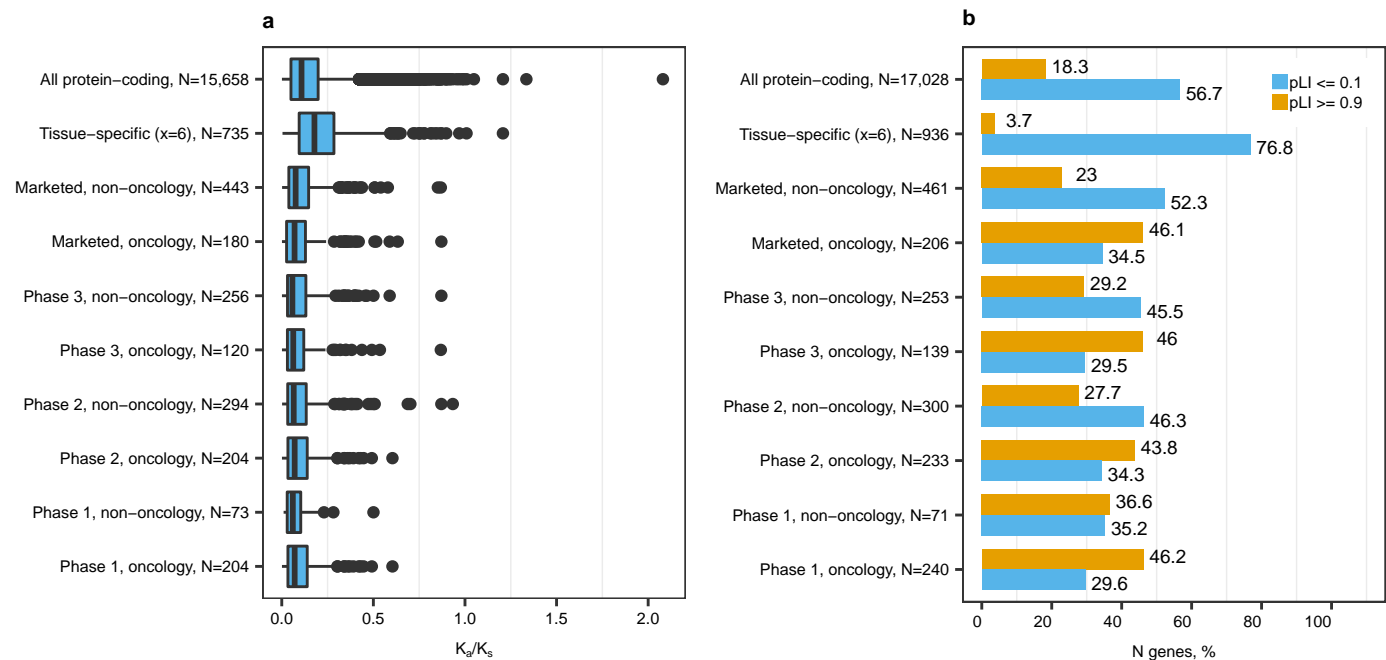
**Fig. 4. Tissue-specific genes were subject to less strong selection pressure compared to all protein-coding genes and drug targets.** **a** Ka/Ks ratios for human-mouse 1:1 orthologs. 1:1 ortholog refers to a human gene with one unique counterpart in mouse as opposed to 1-to-many or many-to-many orthologs that arise from duplication or gene fusion events. **b** Percentages of loss-of-function intolerant (ExAC consortium pLI $\geq$ 0.9) and tolerant (pLI $\leq$ 0.1) genes among tissue-specific genes and drug targets compared to all protein-coding genes. Discrepancies in sample size are due to different numbers of genes mapped to the respective data sets.
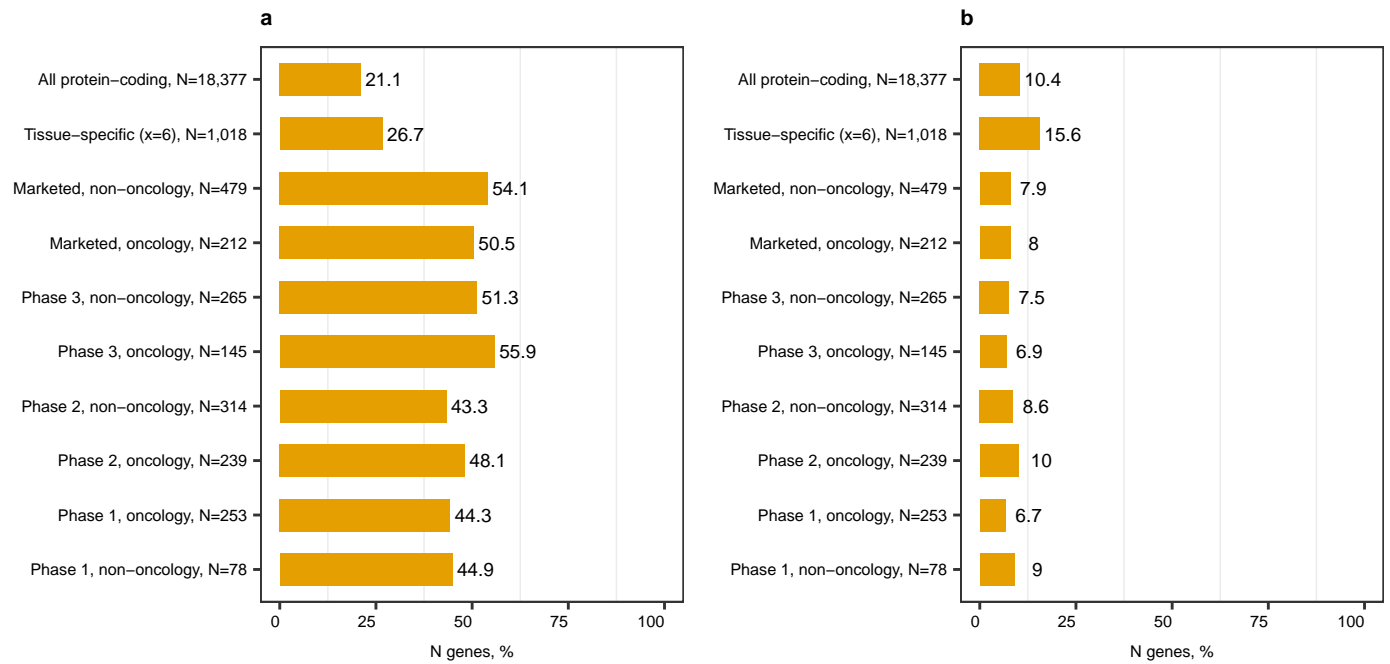
**a**

| Category | Value |
|---|---|
| All protein–coding, N=18,377 | 21.1 |
| Tissue–specific (x=6), N=1,018 | 26.7 |
| Marketed, non–oncology, N=479 | 54.1 |
| Marketed, oncology, N=212 | 50.5 |
| Phase 3, non–oncology, N=265 | 51.3 |
| Phase 3, oncology, N=145 | 55.9 |
| Phase 2, non–oncology, N=314 | 43.3 |
| Phase 2, oncology, N=239 | 48.1 |
| Phase 1, oncology, N=253 | 44.3 |
| Phase 1, non–oncology, N=78 | 44.9 |

N genes, %

**b**

| Category | Value |
|---|---|
| All protein–coding, N=18,377 | 10.4 |
| Tissue–specific (x=6), N=1,018 | 15.6 |
| Marketed, non–oncology, N=479 | 7.9 |
| Marketed, oncology, N=212 | 8 |
| Phase 3, non–oncology, N=265 | 7.5 |
| Phase 3, oncology, N=145 | 6.9 |
| Phase 2, non–oncology, N=314 | 8.6 |
| Phase 2, oncology, N=239 | 10 |
| Phase 1, oncology, N=253 | 6.7 |
| Phase 1, non–oncology, N=78 | 9 |

N genes, %

**Fig. 5. Tissue-specific genes were enriched in disease genes and potential disease genes with gain-of-function mechanism.** The bars show percentages of **a** OMIM and **b** PTVesc genes in each gene category.
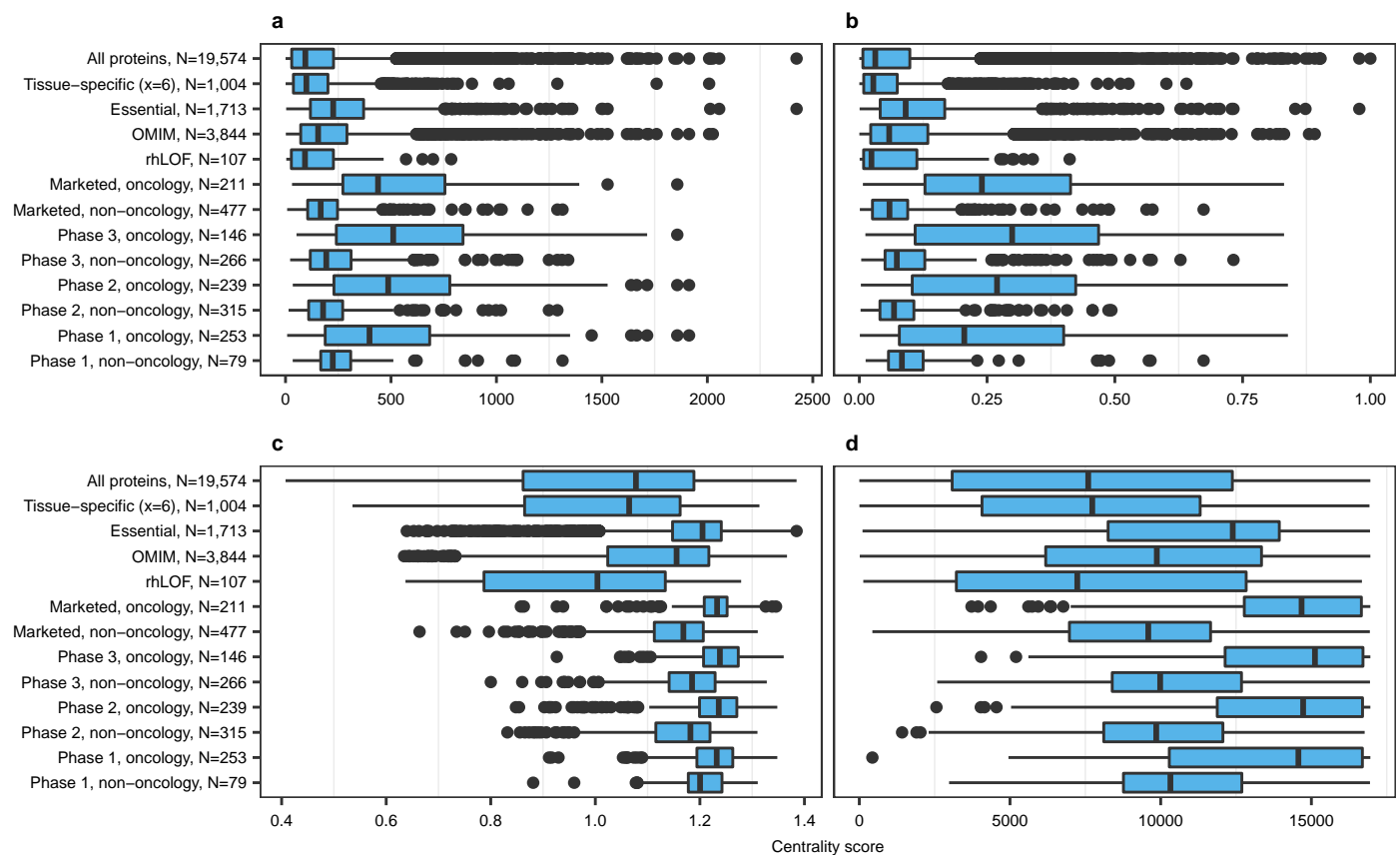
**Fig. 6. Centrality scores in STRING v 10.5. a** Strength **b** Eigenvector centrality **c** Closeness centrality (normalized) **d** Weighted k-shell. Discrepancies in sample size are due to different numbers of genes mapped between data sets.