



Long-read sequencing reveals widespread novel splicing and neojunction-derived neoantigens in nasopharyngeal carcinoma

Yi Shuai, Hualiang Yao, Bo Wang, et al.

Genome Res. published online July 8, 2026

Access the most recent version at doi:[10.1101/gr.281467.125](https://doi.org/10.1101/gr.281467.125)

P<P	Published online July 8, 2026 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Comprehensive immune receptor profiling.
Discover the **DriverMap™ AIR Assay** difference.

LEARN
MORE



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 Long-read sequencing reveals widespread novel splicing and neojunction-derived 2 neoantigens in nasopharyngeal carcinoma

3 Yi Shuai^{1,2,4}, Hualiang Yao^{1,4}, Bo Wang^{2,3}, Grace TY Chung^{2,3}, Xiangeng Wang^{1,4,6}, Ming Zhong^{1,4},
4 Zhongxu Zhu¹, Cheuk Shuen Li⁵, Chi Man Tsang^{2,3}, Kwok-Wai LO^{2,3*}, Xin Wang^{1,4,6*}

5 ¹Department of Surgery, The Chinese University of Hong Kong, Shatin, Hong Kong SAR, China.

6 ²Department of Anatomical and Cellular Pathology, The Chinese University of Hong Kong, Hong
7 Kong SAR, China.

8 ³State Key Laboratory of Translational Oncology, Sir YK Pao Centre for Cancer, The Chinese
9 University of Hong Kong, Hong Kong SAR, China

10 ⁴Li Ka Shing Institute of Health Sciences, The Chinese University of Hong Kong, Shatin, Hong Kong
11 SAR, China

12 ⁵School of Biomedical Sciences, University of Hong Kong, Pokfulam, Hong Kong SAR, China

13 ⁶Shenzhen Research Institute, The Chinese University of Hong Kong, Shenzhen 518000, China

14 * Correspondence: xinwang@cuhk.edu.hk; kwlo@cuhk.edu.hk

15 Abstract

16 The widespread transcriptomic diversity driven by alternative splicing (AS) contributes to all
17 hallmarks of cancer and represents a critical source of neoantigens for personalized immunotherapy.
18 However, unlike other major malignancies, the full repertoire of alternative splicing in nasopharyngeal
19 carcinoma (NPC) remains underexplored. Here, we employ long-read sequencing (LR-seq) to generate
20 a high-resolution, isoform-level transcriptomic atlas from a cohort of 14 NPC tumor samples and four
21 immortalized nasopharyngeal epithelial cell lines. We identify a substantial number of full-length
22 novel transcripts (22,687; ~44.38%), which reveal diverse splicing patterns and previously
23 unannotated splicing events. By integrating short-read RNA-seq data to quantify isoform expression,
24 we discover a subset of novel transcripts that are differentially expressed between tumor samples and

25 immortalized nasopharyngeal epithelial cell lines. Furthermore, LR-seq enables precise identification
26 of chimeric read-through fusion transcripts, such as *CLDN15-FIS1* and *FOXRED2-TXN2*. Finally, we
27 develop a computational framework-**T**umor-**s**pecific **s**plicing **n**eo**a**ntigen **d**etection (TS-SNAD) to
28 predict neoantigens originating from novel exon-exon junctions (neojunctions) in tumor-specific novel
29 transcripts. Using this framework, we identify neojunction-derived neoantigens and experimentally
30 validate the immunogenicity of selected HLA-B*40:01-restricted neoantigens. These neojunction-
31 derived peptides constitute a new class of non-canonical neoantigens with significant potential for
32 developing personalized cancer vaccines for NPC.

33 **Keywords:** nasopharyngeal carcinoma; long-read sequencing; short-read sequencing; alternative
34 splicing; isoform diversity; TS-SNAD; neojunction-derived neoantigens; cancer immunotherapy.

35 Introduction

36 Nasopharyngeal carcinoma (NPC), which originates from the epithelium of the nasopharynx, is
37 endemic in Southeast Asia and Southern China. Distinct from other head and neck cancers, NPC is
38 characterized by its consistent association with Epstein-Barr virus (EBV) infection, unique genomic
39 landscapes and specific histopathological features (Tsang et al. 2020). As early-stage NPC is often
40 asymptomatic, approximately 80% of patients present with locoregionally advanced or metastatic
41 disease at diagnosis. For these patients, chemoradiotherapy is the standard of care, but treatment
42 outcomes remain unsatisfactory for those with locoregionally advanced or metastatic disease.
43 Consequently, novel therapeutic strategies, such as targeted therapies and immunotherapies, are
44 urgently needed (Siak et al. 2023).

45 Alternative splicing (AS) is a key driver of oncogenesis, promoting cancer progression through
46 isoform switches that directly regulate core cancer hallmarks. Most tumors exhibit widespread splicing
47 dysregulation, altering exon-intron architecture, promoter usage, and isoform expression, ultimately
48 disrupting the function of critical tumor suppressors and oncogenes (Huang et al. 2021; Bradley and
49 Anczukow 2023; Maurin et al. 2023; Joglekar et al. 2024; Naro et al. 2024). A comprehensive
50 characterization of transcriptomic diversity is therefore essential for understanding cancer

51 development (Bhattacharya et al. 2023). Conventionally, the detection and quantification of AS
52 events in cancers have relied on short-read RNA sequencing (RNA-seq) (Hong et al. 2020), which is
53 inherently limited by its dependence on algorithmic reconstruction of transcripts from short fragments
54 aligned to a reference genome. As a result, this strategy often yields an incomplete view of the splicing
55 repertoire, as short reads are frequently insufficient to unambiguously resolve complex splice
56 junctions or capture full-length transcripts (Hook and Timp 2023). Recent advances in long-read
57 sequencing (LR-seq) technologies, such as Pacific Biosciences (PacBio) single-molecule real-time
58 (SMRT) sequencing and Oxford Nanopore Technologies (ONT) nanopore sequencing, are poised to
59 overcome these challenges (Monzo et al. 2025; Somalraju et al. 2025). By generating reads that span
60 multiple kilobases, LR-seq enables the direct, accurate identification of full-length alternatively
61 spliced transcripts, allowing for a more comprehensive characterization of the AS landscape in cancer.

62 Several studies have reported that tumor-specific or -associated non-synonymous mutations, gene
63 fusions, transposable elements, and post-translational modifications could generate neoantigens
64 (Burbage et al. 2023; Xie et al. 2023; Hu et al. 2024; Kumar et al. 2024). In practice, neoantigens
65 derived from non-synonymous single-nucleotide variants (SNVs) and gene fusions have been the
66 primary focus of personalized cancer vaccine development in melanoma (Lauss et al. 2020; Ott et al.
67 2020; Borgers et al. 2025), glioma (Platten et al. 2021; Dunn et al. 2022; Zelba et al. 2025), and
68 pancreatic cancer (Sethna et al. 2025). More recently, alternative RNA splicing has emerged as a
69 validated source of neoantigens, including immunogenic cryptic peptides arising from unannotated
70 splice junctions, thereby expanding the antigen repertoire and eliciting antitumor immune responses
71 (Smart et al. 2018; Lu et al. 2021; Merlotti et al. 2023). Despite these novel insights, the non-canonical
72 peptide landscape remains incompletely defined. Previous investigations and computational pipelines
73 for splicing-derived neoantigen discovery have largely relied on short-read RNA-seq (Pan et al. 2023)
74 and reference-based annotations, which require transcript assembly from fragmented reads and may
75 compromise the accurate reconstruction of full-length transcripts and their open reading frames
76 (ORFs).

77 In this study, we employ LR-seq to provide the first systematic view of transcriptomic diversity in
78 NPC, leading to the discovery and characterization of full-length transcripts, including numerous
79 previously undescribed transcripts. By directly resolving full-length transcript structures, LR-seq also
80 reveals widespread alternatively spliced isoforms and read-through transcripts (RTTs). To further
81 validate and quantify the LR-seq-derived transcriptome, we integrate short-read RNA-seq data and
82 identify a tumor-specific subset. Most importantly, we develop **Tumor-specific splicing neoantigen**
83 **detection (TS-SNAD)**, a new computational framework that leverages LR-seq to identify tumor-
84 specific, novel exon-exon junctions (neojunctions) from novel transcripts by directly resolving
85 transcript structures. Our approach minimizes errors introduced by transcript assembly in short-read
86 RNA-seq analyses and enables accurate prediction of neojunction-derived neoantigens. We further
87 employ TS-SNAD to predict neojunction-derived neoantigens and experimentally validate the
88 immunogenicity of selected candidates predicted to bind the high-frequency NPC allele HLA-B*40:01.
89 Our study provides novel insights into the NPC transcriptome and a valuable resource for developing
90 off-the-shelf cancer immunotherapies for NPC patients.

91 Results

92 Long-read sequencing unveils a vast repertoire of novel transcripts and extensive 93 splicing diversity in the NPC transcriptome

94 To comprehensively profile the AS landscape in NPC, we performed long-read RNA sequencing on a
95 cohort of 18 samples, including both non-malignant and malignant samples. This cohort included two
96 NPC cell lines, ten patient-derived xenografts (PDXs), two NPC tumor tissues, and four immortalized
97 nasopharyngeal epithelial (NP) cell lines. All samples were analyzed using SMRT circular consensus
98 sequencing (CCS) on the PacBio Sequel II platform. On average, ~783,000 CCS reads were obtained
99 per library, yielding ~350,000 clustered isoforms after downstream processing (**Table 1 and**
100 **Supplemental Table S1**). For each sample, the numbers of genes and full-length transcripts identified
101 in LR-seq after collapsing were reported (**Supplemental Fig. S1B**). The observed differences across
102 samples were expected and likely reflected variation in RNA quality and biological heterogeneity.

103 After merging transcripts across samples and collapsing redundant isoforms, we identified 51,115
104 high-quality full-length transcripts with a median transcript length of approximately 3.0 kb (**Fig. 1A**
105 **and Fig. 1C**), constructing a comprehensive LR-seq-derived transcriptome for NPC.

106 Transcripts were categorized using SQANTI3 based on their structural alignment to the GENCODE
107 v38 reference transcriptome. This classification revealed 24.10% (12,319) as full splice match (FSM),
108 perfectly matching a known reference transcript; 26.10% (13,341) as incomplete splice match (ISM),
109 which align to a known transcript but lack complete 5' or 3' ends; 22.65% (11,578) as novel in catalog
110 (NIC), containing known splice sites in novel combinations; 21.73% (11,109) as novel not in catalog
111 (NNC), featuring at least one previously unannotated splice site. The remaining 2,768 transcripts were
112 classified as antisense, genic, fusion, or intergenic transcripts (**Fig. 1A**). Novel transcripts (NIC/NNC)
113 constituted a substantial portion of the transcriptome, representing between 17% and 36% of all
114 transcripts per sample (**Fig. 1B**), underscoring the substantial unannotated transcriptional complexity
115 in NPC. Collectively, these transcripts were mapped to 11,807 annotated genes and 2,259 novel
116 genomic loci. Notably, 56.66% (6,690) of the genes in the LR-seq-derived transcriptome expressed
117 two or more transcripts, indicating widespread alternative splicing (**Fig. 1D**).

118 To assess the reliability of the LR-seq-derived transcriptome, we aligned short-read RNA-seq data
119 from matched samples to our LR-seq-derived transcriptome. This validation revealed a clear hierarchy
120 of support [TPM (transcripts per million) ≥ 1]: FSM transcripts showed the highest validation rate
121 (~68-80%); NIC showed intermediate validation rates (~40-60%); NNC and ISM transcripts showed
122 lower support (~30-40%); Antisense, genic, fusion, or intergenic transcripts exhibited the lowest
123 validation rates (~10-25%) (**Fig. 1E**). We next evaluated the distribution of isoforms across samples.
124 We quantified the number of samples supporting each isoform using two criteria: long-read
125 presence/absence calls (binary 0/1 per sample) and short-read RNA-seq expression (counting samples
126 with TPM above the multiple thresholds). FSM transcripts exhibited low heterogeneity, with many
127 isoforms detected across multiple samples. Conversely, novel transcripts exhibited high sample-to-
128 sample variability, with the majority detected in only a few samples, highlighting their greater
129 heterogeneity relative to well-annotated FSM isoforms (**Supplemental Fig. S2A-B**).

130 Comprehensive characterization of novel transcripts with quality evaluation reveals
131 their structural diversity, functional relevance, and potential roles in NPC pathogenesis
132 Next, we performed quality assessments of the identified transcripts, focusing on coding potential,
133 susceptibility to nonsense-mediated mRNA decay (NMD), and splice site characteristics. We observed
134 that the vast majority of FSM and novel transcripts were predicted to be protein-coding (94-97%),
135 whereas antisense, genic, fusion, or intergenic transcripts largely lacked coding potential. NMD
136 prediction revealed that NNC, antisense, genic, fusion, or intergenic transcripts were more likely to
137 contain premature termination codons (PTCs) and non-canonical splice sites (**Fig. 2A**), consistent with
138 increased susceptibility to RNA surveillance and degradation (Supek et al. 2021). Investigating the
139 mechanisms triggering NMD indicated that the use of novel splice sites and intron retention were the
140 primary causes, often generating PTCs that mark transcripts for degradation (**Supplemental Fig. S2C**).
141 For example, two novel *ANO9* transcripts (TCONS_00009816 and TCONS_00009828) were
142 predicted to undergo NMD due to a novel splice site and intron retention, respectively, both
143 introducing PTCs (**Supplemental Fig. S2D**). Furthermore, we observed a positive correlation between
144 the number of exons per gene and the number of novel transcripts produced by that gene, suggesting
145 that genes with higher exon complexity have a greater potential for generating novel splice variants
146 (**Supplemental Fig. S2E**).

147 Comparison with reference transcripts suggested that NNC isoforms often arose from unannotated
148 transcriptional start sites (TSSs) and termination sites (TTSs) located >100 bp from known transcript
149 boundaries (**Supplemental Fig. S2F**). In addition to RNA-seq, we integrated multiple orthogonal
150 datasets to validate novel isoforms identified by LR-seq, including cap analysis of gene expression
151 (CAGE) and 3'-seq peaks from the poly(A) site database. The 5' ends of many novel isoforms
152 overlapped with FANTOM5 CAGE-defined TSSs, and their 3' ends were validated by bona fide TTSs
153 mapped by 3'-seq poly(A) databases (**Supplemental Fig. S2G**). Comparative analysis revealed
154 distinct structural features distinguishing FSM transcripts from novel transcripts. NNC isoforms were
155 typically shorter in total length than FSM and NIC transcripts, yet novel transcripts often possessed
156 longer coding sequences (CDSs) than FSM isoforms. Furthermore, splice junctions in FSM transcripts

157 showed higher read coverage, while novel isoforms, particularly those with novel splice sites, had
158 lower read support (**Fig. 2B**).

159 To evaluate the biological relevance of novel transcripts, we performed pathway enrichment analysis.
160 Genes associated with novel transcripts were significantly enriched in pathways related to RNA
161 processing, nuclear export, and intracellular transport, based on a matched background gene set (**Fig.**
162 **2C**). Notably, several established oncogenes and tumor suppressors, including *FGFR1*, which can be
163 activated by LMP1 to promote glycolysis and invasion in NPC (Lo et al. 2021), *DDX5*, implicated in
164 NAT10-driven immune suppression (Xie et al. 2025), and *MST1R*, which is linked to early-onset NPC
165 and innate immunity (Dai et al. 2016; Cazes et al. 2022) exhibited substantial numbers of novel
166 isoforms, suggesting a crucial role in NPC pathogenesis (**Fig. 2D**). Most novel transcripts (67.8%)
167 shared the same translation start sites as known protein isoforms, indicating that the associated
168 splicing events mainly affected internal exons. In contrast, 32.2% potentially used alternative
169 translation start sites, possibly associated with alternative promoter usage or first exon selection,
170 which may alter the 5' UTR or the exon harboring the canonical start codon (**Fig. 2E**).

171 In addition to novel host transcripts identified in NPC samples, we also detected previously
172 unannotated EBV transcripts, including transcripts with alternative first exon in *LMP2* and exon
173 skipping of exons 2 and 3 in *RPMS1* (**Supplemental Fig. S2H**). To assess the differences between
174 sequencing approaches, we compared the transcriptome derived from LR-seq with the transcriptome
175 assembled from short-read data using StringTie. The StringTie-assembled transcriptome comprised
176 239,952 transcripts, far exceeding the number identified by LR-seq. However, only 14,198 transcripts
177 were commonly identified by both methods, most of which were FSMs. This discrepancy underscored
178 the substantial variation in transcript identification between short-read and long-read sequencing
179 technologies, highlighting the limitations of short-read-based assembly in reconstructing full-length
180 isoforms (**Supplemental Fig. S2I**).

181 Moreover, a key advantage of LR-seq lies in its ability to accurately identify RTTs by capturing full-
182 length transcript structures. This enabled precise delineation of fusion transcripts spanning adjacent
183 genes. These RTTs and their encoded products have emerged as potential diagnostic biomarkers and

184 therapeutic targets in various cancers (Alpert et al. 2020; Hu et al. 2022). In our study, we identified
185 and experimentally validated four such RTTs, including *SFT2D2-TBX19*, *FOXRED2-TXN2*, *ELAC1-*
186 *SMAD4* and *CLDN15-FISI*, using RT-PCR with specific primers targeting the fusion junctions
187 **(Supplemental Fig. S3A-B).**

188 Next, we annotated seven types of alternative splicing events (ASEs) using SUPPA2, including
189 skipping exons (SE), alternative 3' splice sites (A3), alternative 5' splice sites (A5), alternative first
190 exon (AF), alternative last exon (AL), intron retention (RI), and mutually exclusive exons (MX)
191 categories **(Fig. 3A)**, in GENCODE v38 transcriptome and its merged transcriptome incorporating
192 novel transcripts identified by LR-seq. Notably, incorporating novel isoforms increased the number of
193 detected ASEs, particularly AF and SE, while the counts of AL and MX remained largely unchanged
194 **(Fig. 3B)**. Differential ASE analysis between NPC tumor samples and immortalized NP cell lines
195 identified 102 differentially spliced events ($P < 0.05$, $|\Delta\text{PSI}| > 0.1$; **Supplemental Table S2**), with AF
196 (43, 42.16%) and SE (25, 24.51%) being the most prevalent **(Fig. 3C)**. A prime example was *TLNI*,
197 which exhibited a tumor-specific switch involving inclusion of a cancer-associated cassette exon (17b).
198 This isoform, which enhances vinculin binding and cell motility (Gallego-Paez et al. 2023), was
199 significantly upregulated in NPC samples **(Fig. 3C-D)**. We validated this isoform switch by RT-qPCR,
200 confirming that the novel isoform (TCONS_00049374) accounted for a significantly larger fraction of
201 total *TLNI* expression in tumors than the canonical isoforms (ENST00000314888.10) **(Fig. 3E)**.

202 Long-read transcriptomics identifies candidate transmembrane therapeutic targets and 203 tumor-specific transcripts

204 To identify transcripts enriched in NPC tumors, we quantified expression by mapping matched short-
205 read RNA-seq data to our LR-seq-derived transcriptome. Analysis of transcript expression density
206 curves and distribution patterns revealed that most FSMs exhibited relatively higher expression levels
207 compared with novel transcripts **(Fig. 4A-B)**. Using differential transcript expression (DTE) analysis
208 between NPC tumor samples and immortalized NP cell lines, we identified 1,812 significantly
209 dysregulated transcripts ($|\log_2\text{FC}| > 2$, BH-adjusted $P < 0.05$). This included 460 upregulated and 325
210 downregulated FSM transcripts, alongside 241 upregulated and 117 downregulated NIC transcripts,

211 and 211 upregulated and 86 downregulated NNC transcripts (**Supplemental Table S3**). Highly
212 expressed alternatively spliced isoforms in NPC included known isoforms of *CSF2RB*, *LUM*, and
213 *COLIA2*, as well as novel isoforms of *PRRX1*, *LGALS9*, and *CD74* (**Fig. 4C**). Notably,
214 overexpression of *CSF2RB* and *PRRX1* contributes to promoting oncogenic activity and maintaining
215 epithelial-mesenchymal transition (EMT) in human cancers (Du et al. 2021; Charlet et al. 2022). High
216 expression of *LGALS9* (Klibi et al. 2009; Kam et al. 2025) and *CD74* (Chow et al. 2022) has been
217 reported to play important roles in promoting immune evasion and remodeling the tumor
218 microenvironment in NPC. While the well-known transcripts were used in previous functional studies,
219 the identification of alternative and novel transcripts, especially the novel *LGALS9* and *CD74*
220 transcripts, in NPC indicated the potential for new oncogenic activities of these transcripts in NPC
221 tumorigenesis. Because of the potential functional difference between novel and known transcripts,
222 further studies on the NPC-specific transcripts are needed to clarify their key oncogenic activities in
223 tumorigenesis.

224 We next prioritized transcripts predicted to encode plasma membrane transmembrane proteins because
225 these proteins are stably expressed at the cell surface, present extracellularly accessible domains, and
226 therefore represent attractive targets for therapeutic antibodies (Hu et al. 2021; Bardia et al. 2024; Neri
227 et al. 2024; Zaidi et al. 2024; Boixareu et al. 2025). Based on predictions of transmembrane topology
228 and plasma membrane localization, 1,333 known transcripts from 859 genes were predicted to encode
229 cell-surface transmembrane proteins. Additionally, 2,159 novel transcripts from 595 genes were
230 predicted to encode cell-surface proteins containing transmembrane domains (**Supplemental Fig. S4A**
231 **and Supplemental Table S4**). Notably, we identified 294 novel transcripts from 95 genes that were
232 predicted to gain cell-surface localization with transmembrane helices, even though these genes were
233 absent from the TCSA reference surfaceome set (Hu et al. 2021). In contrast, we identified 1,378
234 novel transcripts from 360 genes in the curated TCSA reference surfaceome whose encoded protein
235 products were predicted to lose cell-surface localization relative to the annotated protein products of
236 the same genes (**Supplemental Table S5**). DTE analysis revealed significant upregulation of several
237 cell membrane-associated transcripts in NPC tumor samples (e.g., *IGSF9*, *CXADR*, *TMPRSS2*, *INSR*;

238 **Supplemental Fig. S4B, Supplemental Table S3**). Survival analysis of an NPC public dataset
239 (GSE102349) indicated that high expression of these cell membrane-associated transcripts correlated
240 with poorer progression-free survival (PFS), highlighting their potential clinical relevance
241 (**Supplemental Fig. S4C**).

242 Among the predicted plasma membrane transmembrane isoforms, predominant expression of the
243 NTRK2-T1 protein isoform was reported in our recent publication (Wang et al. 2025). The expression
244 of the *NTRK2-T1* transcript was also confirmed by short-read RNA-sequencing, and the corresponding
245 protein expression was confirmed by western blotting in NPC. By IHC, we demonstrated the
246 membrane staining of NTRK2-T1 protein isoform in the NPC xenografts and primary tumors
247 (**Supplemental Fig. S5A-C**). To further validate the newly identified transmembrane CSF2RB protein
248 isoform TCONS_00033496, we examined its protein expression and subcellular localization in vitro.
249 Western blot analysis confirmed expression of FLAG-tagged TCONS_00033496 in both HEK293T
250 and NP69 cells after transfection. Immunofluorescence staining showed that the FLAG signal was
251 predominantly localized at the cell periphery and colocalized with the plasma membrane marker WGA
252 in both cell lines, whereas no specific anti-FLAG signal was detected in the corresponding negative
253 control cells. These findings indicated that TCONS_00033496 encoded a novel CSF2RB protein
254 isoform that was properly translated and targeted to the plasma membrane (**Supplemental Fig. S5D**).

255 We next assessed the tumor specificity of all AS variants. While the majority of transcripts (34,416,
256 67.33%) were expressed in both NPC tumors and immortalized NP cell lines, and 6,487 (12.69%)
257 were lowly expressed in all samples, a significant fraction (9,349 transcripts, 18.29%) was exclusively
258 detected in NPC tumor samples (**Supplemental Table S6**). In contrast, only 1.69% (863 transcripts)
259 were specific to immortalized NP cell lines (**Fig. 4D**). Tumor-specific novel transcripts constituted a
260 substantial proportion of this group, suggesting they may play a functional role in tumor biology and
261 represent a reservoir of novel biomarkers.

262 To validate these findings, we selected a subset of highly expressed, tumor-specific novel isoforms
263 that showed TPM ≥ 10 in at least one tumor sample and harbored isoform-specific novel junctions
264 with short-read splice junction coverage ≥ 10 . A heatmap of their expression patterns in our RNA-seq

265 cohort confirmed their tumor-specific expression, though substantial inter-tumor heterogeneity was
266 observed (**Fig. 4E**). This tumor specificity was further validated by RT-qPCR (**Fig. 4F**), and the
267 isoform structures were depicted in **Supplemental Fig. S6**. We further validated these results by
268 integrating and analyzing three public NPC RNA-seq datasets (GSE118719, GSE134886, GSE68799),
269 comprising 52 NPC tumors and 11 non-tumor controls after batch-effect correction (**Supplemental**
270 **Fig. S7A**). The expression patterns of the novel transcripts in these independent cohorts were
271 consistent with our in-house data, confirming their specific and elevated expression in NPC tumors
272 (**Supplemental Fig. S7B**).

273 A computational pipeline identifies potential tumor-specific neojunction-derived
274 neoantigens for NPC immunotherapy development

275 To expand the repertoire of immunotherapeutic targets, we investigated neojunctions derived from
276 tumor-specific novel transcripts. These neojunctions represented a source of non-canonical tumor-
277 specific antigens, significantly broadening the landscape of potential targets for cancer therapy.
278 Among all detected splice junctions, the majority (94.1%) were known, while 38,193 (5.90%) were
279 novel (**Fig. 5A**). Most transcripts identified with novel splice junctions contained only a single novel
280 junction, with progressively fewer transcripts carrying multiple novel junctions. Notably, most novel
281 junctions were supported by canonical splice sites (85.3%), while only a minority involved non-
282 canonical splice sites (14.7%), indicating that novel splicing largely conformed to established splice
283 site rules (**Fig. 5B**). Moreover, most novel junctions were derived from novel transcripts, with few
284 originating from intergenic, genic, fusion, or antisense regions (**Fig. 5C**).

285 To systematically identify tumor-specific neoantigens derived from these neojunctions, we developed
286 a customized computational workflow, TS-SNAD. This pipeline identified novel splice junctions from
287 tumor-specific novel transcripts by comparison with reference genome annotations, as well as
288 excluding any splice junctions present in normal tissues from GTEx short-read RNA-seq and long-
289 read RNA-seq data from 88 samples from GTEx tissues (Glinos et al. 2022), thereby reducing the
290 likelihood of selecting peptides with potential expression in vital normal organs and minimizing off-
291 tumor toxicity risk. Tumor specificity was further supported by significantly higher read coverage of

292 these novel splice junctions in NPC tumor samples compared to four immortalized nasopharyngeal
293 epithelial cell lines in our cohort. The pipeline subsequently generated neojunction-spanning 9-mer
294 peptides as candidate neoantigens. The binding affinity of these candidate neoantigens to HLA-I
295 molecules was predicted using NetMHCpan, with an $IC_{50} < 500$ nM and a percentile rank < 0.5 set as
296 selection criteria (**Fig. 5D**). To assess the binding potential of neojunction-derived neoantigens, we
297 compiled HLA-I genotypes from 70 NPC patients from our in-house cohort with whole-genome
298 sequencing data (Bruce et al. 2022) and predicted HLA-I types from 113 Chinese patients with RNA-
299 seq data from the GSE102349 dataset. Despite the high polymorphism of HLA-I, both cohorts showed
300 a consistent presence of high-frequency alleles, including HLA-A*02:07, HLA-A*11:01, HLA-
301 B*46:01, HLA-B*40:01, HLA-C*01:02, and HLA-C*07:02 (**Supplemental Table S7**), likely
302 reflecting the population-specific genetic background.

303 Using TS-SNAD, we identified 944 tumor-specific novel transcripts harboring novel splice junctions
304 with read coverage enriched in tumor samples. From these, we generated a total of 3,326 neojunction-
305 derived 9-mer neoantigens, of which 242 neojunction-derived 9-mer neoantigens matched peptide
306 sequences encoded by other coding regions in the UniProtKB/Swiss-Prot human proteome. Of the
307 3,326 neojunction-derived 9-mer candidates, 533 candidate neoantigens were predicted to bind
308 patient-specific HLA-I alleles with high affinity, suggesting a high likelihood of being presented (**Fig.**
309 **5E**). Notably, most of these neoantigens were derived from unique tumor-specific novel transcripts,
310 while 70 were produced by multiple isoforms (**Fig. 5F**). Analyzing the diversity of HLA-I interactions,
311 we found that most neoantigens could bind multiple HLA-I alleles (**Fig. 5G**), suggesting their
312 potential to be presented across a broader patient population with diverse immune contexts, thereby
313 enhancing their therapeutic relevance.

314 Neojunction-derived neoantigens effectively elicit immunogenic responses

315 To evaluate the potential of neojunction-derived neoantigens as targets for broadly applicable
316 immunotherapies in NPC, we selected 34 neojunction-derived peptides with high predicted binding
317 affinity for HLA-B*40:01, a high-frequency allele present in approximately 29.20% of NPC patients,
318 and assessed their immunogenicity. Using peripheral blood mononuclear cells (PBMCs) from three

319 HLA-B*40:01-positive healthy donors, we measured interferon- γ (IFN- γ) responses following
320 stimulation with peptide-MHC (pMHC) complexes. Six of the 34 predicted neoantigen-derived
321 neoantigens induced robust and reproducible IFN- γ responses in all three donors (**Fig. 6A,**
322 **Supplemental Fig. S8A and Supplemental Table S8**). The absence of response in DMSO controls
323 and the robust activation induced by a CD3-specific antibody confirmed the specificity and reliability
324 of the assay. The most immunogenic peptide, REPFVQAKL, was derived from a novel *TYK2*
325 transcript (TCONS_00025200) generated through skipping of exons 7-9. We observed
326 considerable inter-donor variability in the magnitude of IFN- γ responses, with donor 3
327 showing consistently stronger responses across all peptides, likely reflecting differences in
328 individual immunological history.

329 Of note, each immunogenic neoantigen originated from a distinct tumor-specific novel
330 transcript in our cohort, with the majority of underlying novel splice junctions resulting from
331 exon skipping or cassette exon inclusion (**Fig. 6B and Supplemental Fig. S8B(a)**). The novel
332 transcripts giving rise to these immunogenic neoantigens exhibited multiple features
333 supporting their translational potential, including canonical splice sites, methionine-initiated
334 ORFs, and no predicted susceptibility to nonsense-mediated mRNA decay. Analysis of public
335 NPC datasets confirmed that these novel isoforms exhibited high tumor-specific expression
336 (**Supplemental Fig. S8B(b)**), minimizing potential risks of off-target reactivity against
337 control tissues. Neoantigen-derived neoantigens largely recapitulated the canonical HLA-
338 B*40:01 binding motif, with a dominant glutamate at P2 and hydrophobic residues at the C
339 terminus (**Supplemental Fig. S8C**).

340 As each neoantigen was found to be present in only 10-20% of NPC patients, these results support the
341 development of personalized multi-epitope vaccines tailored to individual splicing profiles to
342 maximize immune targeting and overcome immune evasion. Furthermore, we confirmed the presence
343 of neoantigen-derived unique peptides in NPC samples via mass spectrometry (**Fig. 7A**). It is

344 noteworthy that junction-spanning peptides are typically underrepresented in mass spectrometry
345 datasets due to the cleavage specificity of trypsin (Kahles et al. 2018), underscoring the significance of
346 these findings. Additionally, to assess the protein-coding potential of these novel RNA isoforms that
347 generated validated immunogenic neoantigens, we overexpressed their ORFs fused to a FLAG tag in
348 both HEK293T and NP69 cells. Western blot analysis showed that all selected novel ORFs generated
349 detectable proteins in vitro (**Fig. 7B**). Collectively, our results demonstrated that neojunction-derived
350 neoantigens constituted a rich source of immunogenic targets, highlighting their strong potential for
351 NPC immunotherapy development.

352 Discussion

353 In this study, we leveraged LR-seq to construct a comprehensive, isoform-level transcriptomic atlas of
354 NPC, revealing widespread transcript diversity and novel splicing events that were largely
355 undetectable by short-read sequencing. We validated the reliability of these full-length transcripts and
356 demonstrated that previously unannotated isoforms, particularly those classified as NIC, NNC, and
357 read-through chimeric transcripts, substantially contributed to the transcriptomic complexity of NPC.
358 Integrated analysis with RNA-seq data revealed that many of these novel isoforms exhibited tumor-
359 specific expression patterns, some of which were strongly associated with clinical outcomes,
360 underscoring their potential as biomarkers and therapeutic targets. We selected a subset of tumor-
361 specific novel transcripts for RT-qPCR validation using isoform-specific primers and observed a
362 relatively high overall validation rate. The inability to validate a small subset of candidates was likely
363 due to two practical factors. First, for genes with multiple isoforms sharing extensive exonic regions,
364 isoform-level abundance estimates inferred from short-read RNA-seq using the LR-seq-derived
365 transcriptome as a reference may be less accurate, which could have affected the reliable assessment
366 of tumor specificity. Second, successful qPCR primer design was not always feasible, because it had
367 to satisfy not only standard primer-design criteria, including appropriate GC content, minimal primer-
368 dimer or hairpin formation, high specificity, and suitable amplicon length, but also the requirement to
369 uniquely distinguish the isoform of interest from other highly similar transcripts. Therefore, failure to
370 validate a small subset of candidates by qPCR did not necessarily indicate that these transcripts were

371 false positives, but more likely reflected technical limitations in isoform quantification and assay
372 design.

373 In contrast to conventional neoantigens derived from gene mutations, we identified a rich source of
374 tumor-specific antigens originating from neojunctions within novel transcripts, which were efficiently
375 detected by LR-seq. These splicing alterations often disrupt canonical open reading frames, giving rise
376 to novel protein products. To systematically harness this potential, we developed the TS-SNAD
377 pipeline. TS-SNAD leverages long-read-resolved full-length transcript structures to more accurately
378 identify novel isoforms and their associated novel splice junctions, providing a more reliable
379 assessment of complex splicing events and their coding consequences than de novo assembly based on
380 short-read RNA-seq or reference-guided transcript reconstruction. In addition, TS-SNAD prioritizes
381 tumor-specific novel junctions by integrating group-wise junction coverage analysis and dual filtering
382 against GTEx short-read and long-read normal tissue exon-exon junctions. Finally, it generates
383 peptides that truly span novel junctions and further evaluates their HLA-binding potential, thereby
384 providing an integrated framework for identifying splicing-derived neoantigens. Using TS-SNAD, we
385 initially identified 3,326 tumor-specific neojunction-derived candidate neoantigens, 533 of which were
386 predicted to bind with high affinity to specific NPC HLA-I alleles and were not present in the
387 reference proteome. TS-SNAD thus provides an effective framework for prioritizing immunogenic
388 neojunction-derived neoantigens from tumor-specific novel transcripts, significantly expanding the
389 repertoire of potential targets for immunotherapy.

390 To advance the development of broadly applicable therapies, we summarized high-frequency HLA-I
391 alleles within the East Asian NPC population, including HLA-A*02:07, HLA-A*11:01, HLA-
392 B*46:01, HLA-B*40:01, HLA-C*01:02, and HLA-C*07:02. This stable HLA landscape provided a
393 foundation for designing personalized yet population-compatible immunotherapeutic strategies.
394 Focusing on HLA-B*40:01, we selected 34 predicted neoantigens for experimental validation. Ex vivo
395 IFN- γ ELISpot assays confirmed that six of these elicited robust immune responses across multiple
396 donors. These results provide further evidence that TS-SNAD can effectively identify immunogenic
397 neojunction-derived neopeptides.

398 Despite these encouraging results, several limitations should be acknowledged. First, the control group
399 comprised immortalized nasopharyngeal epithelial cell lines rather than true normal nasopharyngeal
400 epithelial tissues. Though these cell lines provide a practical epithelial control, they may not fully
401 recapitulate the molecular features of native normal epithelium. In addition, publicly available bulk
402 RNA-seq datasets from normal nasopharyngeal biopsies or adjacent non-malignant tissues are
403 suboptimal for validating the differentially spliced events identified here, because these specimens
404 contain very few nasopharyngeal epithelial cells and are dominated by stromal and immune cells. The
405 lack of RNA-seq data from microdissected normal nasopharyngeal epithelial cells therefore limits
406 direct validation of tumor-associated splicing changes in primary normal epithelial counterparts.
407 Future studies using microdissected normal epithelium, epithelial-enriched samples, single-cell or
408 long-read transcriptomic approaches will be important to further confirm the specificity of these
409 splicing events. Second, we partially validated their immunogenicity using *ex vivo* IFN- γ ELISpot
410 assays. Nevertheless, further *in vitro* and *in vivo* studies will be required to confirm antigen
411 presentation and functional relevance in the tumor context. Third, while LR-seq accurately captures
412 full-length isoforms, its lower throughput and sequencing depth relative to short-read RNA-seq may
413 limit the detection of low-abundance transcripts. Finally, as our HLA frequency and immunogenicity
414 data were primarily derived from Chinese cohorts, extrapolation to other ethnicities requires validation
415 in more diverse populations.

416 In conclusion, our work highlights the critical value of integrating LR-seq into cancer transcriptomics
417 research, not only for discovering novel isoforms and splicing events but also for unveiling a hidden
418 landscape of functionally and clinically relevant neoantigens. The TS-SNAD pipeline establishes a
419 foundation for rational neoantigen selection and paves the way for developing off-the-shelf or
420 personalized cancer vaccines for NPC. Future directions include immunopeptidomics validation of
421 predicted neoantigens, functional assessment of T cell responses, and integration with single-cell
422 transcriptomics and proteogenomics to achieve a more comprehensive immunological understanding.

423 **Methods**

424 **Long-read library preparation and sequencing**

425 A total of 18 samples were collected and subjected to long-read sequencing, including four
426 immortalized nasopharyngeal cell lines (NP361, NP460, NP550, and NP69) as controls, and 14 NPC
427 samples comprising two cell lines (NPC43 and NPC53), 10 patient-derived xenografts (X2117, X17,
428 C15, X111, X113, X666, X76, X47, X23, and X32), and two tumor tissues (NPC-T9 and NPC-T64).
429 This sample set enabled transcript discovery across diverse biological sources. Following RNA
430 extraction, full-length cDNA was synthesized from poly(A)-tailed RNA using the SMARTer PCR
431 cDNA Synthesis Kit (Clontech/Takara; Cat. #634926). The cDNA was PCR-amplified to yield 1 to 2
432 μg of product and purified using AMPure XP magnetic beads (Beckman Coulter; A63881). Size
433 selection was performed using the Sage Science BluePippin system to remove short cDNA fragments
434 and reduce preferential loading and sequencing of shorter molecules. SMRTbell adapters were ligated
435 to cDNA ends, and libraries were purified using magnetic beads with the SMRTbell Template Prep
436 Kit (Pacific Biosciences), followed by sequencing on PacBio Sequel II platform.

437 **RNA sequencing**

438 RNA was extracted using the RNeasy Mini Prep kit (Qiagen; Cat. #74106), quantified on a Qubit 2.0
439 fluorometer (Thermo Fisher Scientific), and assessed for quality using a 2100 Bioanalyzer (Agilent)
440 with the RNA 6000 Pico kit. cDNA libraries were prepared using Illumina TruSeq Stranded mRNA
441 Library Prep Kit. The final libraries were sequenced on an Illumina NovaSeq 6000 platform to
442 generate 150-bp paired-end reads.

443 **Quantitative PCR with reverse transcription**

444 Reverse transcription was performed with a High-Capacity RT kit (Applied Biosystems, Thermo
445 Fisher Scientific; Cat#4368813). The resulting cDNA was diluted 1:3 using nuclease-free water, and 2
446 μL was used per quantitative PCR (qPCR) reaction. qPCR with reverse transcription was performed
447 using the POWER SYBR Green Master Mix (Applied Biosystems, Thermo Fisher Scientific;
448 Cat#4367659). All samples were analyzed in technical triplicates using the QuantStudio 5 (Thermo

449 Fisher Scientific), and all gene expression data were normalized to human *GAPDH*. Primer sequences
450 were listed in **Supplemental Table S9**.

451 Long-read sequencing data analysis

452 Raw PacBio subreads were processed using the Iso-Seq3 pipeline (version 3.8.0)
453 (<https://github.com/PacificBiosciences/IsoSeq/>) to generate high-quality, full-length transcript
454 sequences for downstream analysis. Briefly, raw subreads were processed to generate CCS reads with
455 a minimum predicted accuracy of 0.99. Full-length reads were generated by removing the 5' and 3'
456 cDNA primers and performing demultiplexing using the 'lima' module with default parameters.
457 Artificial concatemers and poly(A) tails were removed using 'isoseq3 refine' to produce full-length,
458 non-chimeric (FLNC) reads. These FLNC reads were clustered into high-quality transcripts using
459 'isoseq3 cluster'. The filtered transcripts were aligned to the human reference genome (GRCh38/hg38)
460 using pbmm2 (version 1.10.0) (<https://github.com/PacificBiosciences/pbmm2>). Redundant isoforms
461 were collapsed within each sample using the collapse_isoforms_by_sam.py script from the
462 cDNA_Cupcake ToFU toolkit (https://github.com/Magdoll/cDNA_Cupcake). Finally, transcripts from
463 all samples were merged into a unified, non-redundant transcriptome using gffcompare (version
464 0.11.2) (Pertea and Pertea 2020) for subsequent analysis (see **Supplemental Fig. S1A**).

465 Isoform annotation and quality control

466 Full-length isoforms were annotated and subjected to quality control using SQANTI3 (version 5.2)
467 (Pardo-Palacios et al. 2024) with GENCODE v38 as the reference annotation. SQANTI3 was used to
468 annotate isoforms and generate quality-control metrics, including CAGE peaks (Noguchi et al. 2017),
469 poly(A) sites (Moon et al. 2025), and NMD predictions, to help filter potential artifacts and ensure
470 high-quality transcript annotations. Isoform validity was further supported by short-read RNA-seq data.
471 The coverage of isoform splice junctions was calculated using the junction file (commonly referred to
472 as SJ.out.tab) generated from RNA-seq data. ORFs were predicted with ORFfinder
473 (<https://github.com/Chokyotager/ORFFinder>), retaining only those encoding ≥ 100 amino acids.

474 De novo transcriptome assembly

475 Sorted BAM files (hg38-aligned) were assembled *de novo* using StringTie (version 2.1.1) (Pertea et al.
476 2015), with parameters: -m 200 -f 0.05 -c 1.5 -p 10. Individual assemblies from the NPC samples were
477 merged with StringTie --merge with default parameters. Both the *de novo* and the LR-seq-derived
478 transcriptomes were compared to GENCODE v38 using gffcompare -G.

479 Short-read RNA sequencing analysis

480 The quality control of short-read RNA-seq data was performed using FastQC (version 0.11.9). The
481 low-quality reads and adapters were filtered and trimmed by using cutadapt (version 4.4) (Kechin et al.
482 2017). Clean reads were mapped to the human reference genome hg38 using STAR (version 2.7.10b)
483 (Dobin et al. 2013). Specifically, a two-pass strategy was used to extract splice junctions from each
484 sample. The TPM abundances of all isoforms were calculated using RSEM (version 1.3.1) (Li and
485 Dewey 2011) using the LR-seq-derived transcriptome as reference. Public RNA-seq data from
486 different cohorts were integrated, followed by batch effect removal using ComBat from the R *sva*
487 package.

488 Alternative splicing event analysis

489 We used SUPPA2 (version 2.3) (Trincado et al. 2018) to calculate seven types of alternative splicing
490 events including A3, A5, AF, AL, SE, RI, and MX. To incorporate novel transcripts, a unified GTF
491 file (merged from GENCODE v38 and long-read-derived novel isoforms) was used with
492 the generateEvents command (-f ioe). Differential splicing analysis was performed with SUPPA2
493 diffSplice, with significance thresholds set at $|\Delta\text{PSI}| > 0.1$ and $P < 0.05$. To improve accuracy, splice
494 sites were held to stringent boundaries, while transcription start/end sites (AF/AL events) allowed
495 flexible boundaries (± 48 nt).

496 Pathway enrichment analysis

497 Pathways significantly enriched (false discovery rate [FDR] < 0.05) among genes with novel isoforms
498 detected by LR-seq were identified using the MSigDB C2, C5, and Hallmark collections (Liberzon et
499 al. 2015). Enrichment analyses were performed using the R package clusterProfiler (version 3.16.1)

500 (Wu et al. 2021). A matched background gene set was used for all enrichment tests and was defined as
 501 the intersection of (i) genes detectable in our long-read data and (ii) genes that could be mapped to the
 502 selected MSigDB collections. This background ensured that enrichment was evaluated relative to
 503 genes testable in our experiment and reduced potential bias.

504 To control for potential bias arising from gene expression levels, we performed an expression-matched
 505 permutation analysis. For each gene, expression level was defined as the median TPM across samples
 506 within the expressed gene universe. Genes were stratified into quantile-based expression bins (10 bins
 507 by default). For the observed gene set generating novel isoforms, we recorded the number of genes in
 508 each expression bin and generated matched random gene sets with the same expression-level
 509 distribution and repeated this procedure across 200 permutations. Pathway enrichment analysis was
 510 then performed for each randomized gene set. For each pathway significantly enriched in the genes
 511 with novel isoforms, an empirical p value (p_{emp}) was calculated as the fraction of expression-matched
 512 permuted gene sets whose p value (p_{perm}) was less than or equal to the observed p value (p_{obs})
 513 obtained for the gene set with novel isoforms:

$$p_{\text{emp}} = \frac{1 + \sum_{b=1}^B 1(p_{\text{perm}}^{(b)} \leq p_{\text{obs}})}{B + 1}$$

514 where B denotes the number of permutations. This analysis provided a direct test of whether the
 515 pathway enrichments observed in genes generating novel isoforms exceeded what would be expected
 516 based solely on their expression-level distribution. Pathways with empirical $p_{\text{emp}} < 0.05$ were
 517 considered significant.

518 Differential transcript expression analysis

519 RSEM expected read counts were used as input for downstream differential expression analysis.
 520 Differential transcript expression between NPC tumor samples and immortalized nasopharyngeal cell
 521 lines was assessed in DESeq2 (v1.38.3) (Love et al. 2014) using the Wald test under a negative
 522 binomial generalized linear model, with Benjamini-Hochberg correction for multiple testing.

523 Transcripts with $|\log_2FC| > 2$ and adjusted $P < 0.05$ were considered significantly differentially
524 expressed.

525 Identification of transcripts encoding putative cell-surface transmembrane proteins

526 An isoform was classified as encoding a putative cell-surface transmembrane protein if its ORF
527 contained at least one transmembrane helix predicted by DeepTMHMM (Hallgren et al. 2022) and
528 was predicted by DeepLoc (Odum et al. 2024) to localize to the cell membrane. For external validation,
529 we used the TCSA Cancer Surfaceome Atlas as a gene-level reference set of genes encoding cell-
530 surface proteins. TCSA integrates nine independent resources to define a curated cancer surfaceome.
531 Because currently available databases generally lack isoform-resolved subcellular localization
532 annotations, we could not reliably assign cell-surface localization to a single canonical isoform for
533 every gene. Therefore, the gain or loss of predicted cell-surface localization was assessed at the
534 isoform level relative to the annotated protein products of the corresponding gene.

535 Identification of splicing-derived neoantigens by TS-SNAD

536 Tumor specificity was assessed by integrating long-read and short-read RNA-seq data, as the
537 sequencing depth of our long-read dataset was insufficient for robust transcript-level abundance
538 estimation. Long-read RNA-seq was used primarily to construct a LR-seq-derived transcriptome by
539 resolving full-length transcript structures and identifying both annotated and previously unannotated
540 isoforms, whereas short-read RNA-seq was used to quantify transcript expression based on the LR-
541 seq-derived transcriptome. Tumor-specific transcripts were defined as transcripts with $TPM \geq 1$ in at
542 least one tumor sample and not detected in any of four immortalized NP cell lines in our cohort. Splice
543 junction coverage was quantified using SJ.out.tab files generated by STAR alignment of RNA-seq
544 data to the hg38 genome, with the LR-seq-derived transcriptome annotation provided as a reference.
545 For neojunction prioritization, all novel splice junctions derived from tumor-specific novel transcripts
546 were subjected to further filtering. Candidate neojunctions were required to meet the following criteria:
547 (i) the average junction read coverage in tumor samples was higher than that in immortalized non-
548 malignant nasopharyngeal epithelial cell lines; (ii) < 5 reads supporting the novel splice junctions in
549 each immortalized NP cell line; and (iii) < 5 reads supporting the novel splice junctions in each of 11

550 non-malignant nasopharyngeal tissue samples from three public NPC RNA-seq datasets (GSE118719,
551 GSE134886, and GSE68799); (iv) the novel splice junctions were required to be absent from all GTEx
552 datasets, including short-read RNA-seq from normal tissues and Oxford Nanopore long-read data from
553 88 GTEx tissues and cell lines.

554 For neoantigen prediction, 9-mer peptides spanning neojunctions were extracted. NetMHCpan
555 (version 4.0) (Jurtz et al. 2017) was then used with default parameters to predict HLA binding affinity,
556 using neojunction-derived 9-mer peptides and patient-specific HLA genotypes as input. Only
557 candidate neoantigen peptides with strong predicted HLA-I binding were retained, defined as those
558 with a binding affinity percentile rank < 0.5 and predicted $IC_{50} < 500$ nM.

559 HLA-I haplotype prediction

560 HLA-I genotypes for 70 patients from our local cohort were determined from our previous whole-
561 genome sequencing data and subsequently used for neoantigen binding affinity prediction. For public
562 datasets, HLA-I genotypes were predicted from RNA-seq FASTQ files using OptiType (version 1.3.3)
563 (Szolek et al. 2014). We selected OptiType based on its demonstrated high accuracy (over 99%) in
564 predicting HLA-I alleles from RNA-seq data.

565 IFN- γ enzyme-linked immunospot (ELISpot) assay

566 The immunogenicity of each neoantigen peptide was evaluated using an IFN- γ ELISpot assay. Briefly,
567 PBMCs from donors carrying the HLA-B*40:01 allele were isolated from whole blood using Ficoll-
568 Paque (Cytiva, 10294426) density gradient centrifugation. The collected PBMCs were then
569 resuspended in RPMI-1640 (Gibco, 25367098) complete medium supplemented with 10% fetal bovine
570 serum (Gibco, 24341305) and incubated at 37°C in a humidified incubator with 5% CO $_2$ to allow cell
571 recovery and stabilization. Subsequently, 3×10^5 PBMCs in 100 μ L were seeded into each well of a
572 96-well ELISpot plate (Mabtech, 3420-4AST-2), followed by the addition of 2 μ L of peptide
573 stimulator (final concentration 20 μ g/mL). Anti-CD3 monoclonal antibody was used as a positive
574 control, and medium containing an equivalent concentration of DMSO served as a negative control.
575 Plates were then incubated at 37°C in a humidified 5% CO $_2$ incubator for at least 20 hours without
576 disturbance. Measures were taken to prevent evaporation during incubation. After incubation, cells

577 were removed, and the plates were washed five times with PBS. A biotinylated anti-IFN- γ detection
578 antibody (1:1000 dilution) was added to each well and incubated for 2 hours at 37°C. After another
579 round of washing, streptavidin-ALP (1:1000) was added and incubated for 1 hour at room temperature.
580 Finally, BCIP/NBT substrate solution was applied to develop colorimetric spots representing IFN- γ -
581 secreting cells. Spot-forming units (SFUs) were quantified using an automated ELISpot reader (AID
582 iSpot Spectrum). Responses were considered positive if spot counts exceeded background by at least
583 threefold.

584 Cell cultures, plasmids and transfection

585 NP69 cells were cultured in Keratinocyte SFM (1X) Kit (Gibco, 17005042), supplemented with 1%
586 penicillin-streptomycin (10,000 U/mL, Gibco, 15140122). HEK293T cells were maintained in
587 DMEM (high glucose, Gibco, 11965118) supplemented with 10% FBS (Sigma-Aldrich, F2442-
588 500ML) and 1% penicillin-streptomycin. All cell lines were cultured at 37°C with 5% CO₂ and were
589 routinely tested and confirmed to be free of mycoplasma contamination.

590 ORFs of the novel transcripts, followed by a 3 \times FLAG tag, were cloned into the pcDNA3.1(+) vector,
591 and all plasmid constructs were verified by sequencing. HEK293T cells were transiently transfected
592 using Lipofectamine 3000 reagent (Invitrogen, L3000001) according to the manufacturer's
593 instructions. NP69 cells were transfected with FuGENE HD (Promega, E2311) according to the
594 manufacturer's instructions.

595 Reagents and antibodies for immunoblotting

596 The following antibodies were obtained from the indicated sources and used for immunoblotting: anti-
597 FLAG M2 mouse mAb (Sigma-Aldrich, F1804) and anti- β -actin rabbit antibody (ABclonal AC038).
598 Secondary antibodies were obtained from the following sources: HRP-conjugated goat anti-mouse IgG
599 (H+L) (Invitrogen, 31430) and HRP-conjugated goat anti-rabbit IgG (H+L) (ABclonal, AS014).

600 Immunofluorescence

601 Cells grown on coverslips were fixed with 4% paraformaldehyde for 10 min at room temperature and
602 washed with PBS. Fixed, non-permeabilized cells were then incubated with wheat germ agglutinin-

603 AF488 (MCE, HY-NP163A) in PBS for 10 min at room temperature to label the plasma membrane.
604 Following WGA staining, cells were washed with PBS and subsequently permeabilized with 0.2%
605 Triton X-100 in PBS for 5 min. Cells were then blocked with 5% BSA in PBS for 60 min at room
606 temperature. Cells were incubated with anti-FLAG M2 mouse monoclonal antibody (Sigma-Aldrich,
607 F1804) overnight, followed by incubation with goat anti-mouse IgG (H+L) secondary antibody
608 (Invitrogen, A11004) for 60 min. Nuclei were stained with DAPI. Coverslips were mounted using
609 antifade mounting medium, and images were captured using a Leica STELLARIS 8 FALCON
610 confocal microscope.

611 Proteomic profiling and peptide database search

612 Public NPC mass spectrometry data were obtained from the iProX partner repository (Accession No.
613 IPX0001265000). Peptide identification was performed using MS-GF+ (version 2018.10.15) (Kim
614 and Pevzner 2014) against a custom sequence database comprising a total of 42,613 merged ORFs,
615 including 22,017 ORFs derived from long-read novel transcripts and 20,596 reviewed human protein
616 sequences from UniProtKB/Swiss-Prot database. The following parameters were set for database
617 searching: Carbamidomethyl (C), iTRAQ4plex (N-term), and iTRAQ4plex (K) were specified as fixed
618 modifications; Oxidation (M), Deamidated (NQ), Acetyl (K), and Methyl (K) were specified as
619 variable modifications. The precursor mass tolerance was set to 20 ppm, and the product ion tolerance
620 for MS/MS was 0.05 Da. Trypsin was specified as the digestion enzyme, allowing up to two missed
621 cleavages. mzID profiles identified from the search engine were then pooled using the R/Bioconductor
622 package MSnbase, and peptide-to-spectrum matches (PSMs) satisfying both spectra and peptide false
623 discovery rate cutoffs <1% were retained for further analysis. To validate neojunction-derived
624 neoantigens by mass spectrometry, only unique 9-mer peptides with unambiguous spectral matches
625 were considered.

626 Survival analysis

627 Survival analysis was performed to evaluate the prognostic associations of transcript expression levels
628 in NPC patients. Patients were stratified into high- and low-expression groups based on the median

629 transcript levels. Kaplan-Meier survival curves were generated and compared using the log-rank test
630 with the survival (version 3.3) and survminer (version 0.4.9) R packages.

631 Data Sets

632 Public RNA-seq datasets from NPC cohorts were obtained from the Gene Expression Omnibus (GEO)
633 with accession numbers GSE118719, GSE134886, GSE68799 and GSE102349. RNA-seq splice
634 junction data from normal human tissues were retrieved from the GTEx project
635 (https://www.gtexportal.org/home/downloads/adult-gtex/bulk_tissue_expression). Long-read RNA-
636 seq data for 88 GTEx tissue and cell lines, based on the Oxford Nanopore Technologies platform,
637 were also retrieved from the GTEx project ([https://www.gtexportal.org/home/downloads/adult-
638 gtex/long_read_data](https://www.gtexportal.org/home/downloads/adult-gtex/long_read_data)). Public mass spectrometry proteomic data of NPC were downloaded from the
639 iProX partner repository (number IPX0001265000).

640 Ethics Declarations

641 For the two NPC tumor samples, written patient consents were obtained from the patients according to
642 institutional clinical research approval. The study protocol was approved by the Joint Chinese
643 University of Hong Kong-New Territories East Cluster Clinical Research Ethics Committee at the
644 Chinese University of Hong Kong, Hong Kong SAR. The collection and use of PBMCs from healthy
645 donors were approved under the protocol “Collection of peripheral blood mononuclear cells for new
646 drug research and development” (protocol No. LP202210). Written informed consent was obtained
647 from all PBMC donors before sample collection.

648 Data Access

649 Long-read RNA sequencing data generated in this study have been submitted to the NCBI BioProject
650 database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA1333213.

651 Code availability

652 All code required to reproduce the data processing and downstream analyses in this study has been
653 provided as Supplemental Code files and archived in a publicly accessible GitHub repository. Within

654 this repository, the analysis scripts are available at [https://github.com/CityUHK-CompBio/NPC-LR-](https://github.com/CityUHK-CompBio/NPC-LR-Seq/tree/main/script)
655 [Seq/tree/main/script](https://github.com/CityUHK-CompBio/NPC-LR-Seq/tree/main/TS-SNAD), whereas the TS-SNAD pipeline is available at [https://github.com/CityUHK-](https://github.com/CityUHK-CompBio/NPC-LR-Seq/tree/main/TS-SNAD)
656 [CompBio/NPC-LR-Seq/tree/main/TS-SNAD](https://github.com/CityUHK-CompBio/NPC-LR-Seq/tree/main/TS-SNAD). The TS-SNAD directory includes the full pipeline
657 implementation, example input and output files, parameter descriptions, and a step-by-step usage
658 guide.

659 Competing interest statement

660 The authors have no conflicts of interest to declare.

661 Acknowledgments

662 This work was funded by grants from Shenzhen Medical Research Funds (C2303002, X.W.), a startup
663 grant (4937084, X.W.), a direct grant (2024.175, X.W.), grants by the Faculty Postdoctoral Fellowship
664 Scheme (FPFS/24-25/053, FPFS/23-24/061C, FPFS/23-24/060, X.W.), Research Committee - Group
665 Research Scheme 2022-23 (WW/rc/grs2223/0560/23en, X.W.), and Faculty Innovation Award
666 (FIA2020/A/01, C.M.T.) from the Chinese University of Hong Kong, and grants from the Research
667 Grants Council (AoE/M-401/20 to X.W. & K.W.L.; R4007-23, C4024-22GF, 14104223, 11103921,
668 and 14111522 to X.W.; 14101721 to K.W.L.; 14124925, 14116124, 14113620, 14114523, and
669 24114922 to C.M.T.), The Innovation and Technology Fund (MRP/036/21X, K.W.L.), and Health and
670 Medical Research Fund (08192166 to X.W.; 08191046 to K.W.L; 09203176 to C.M.T.). This work
671 was also partially sponsored by Shenzhen Bay Scholars Program and Jiangxi Overseas High-Level
672 Talent Project (20232BCJ25029) awarded to X.W. The authors gratefully acknowledge the valuable
673 contributions of all members of Prof. Xin Wang's and Prof. K.W. Lo's laboratory, who provided
674 insightful feedback and shared their expertise in conducting the experiments and performing the
675 bioinformatics computations.

676 Author contributions

677 X.W., Y.S., and K.W.L. conceived the project and designed the research. Y.S. conducted the
678 bioinformatics analyses. H.L.Y., B.W., G.T.Y.C., M.Z. and Y.S. performed the experimental

679 validation. Y.S. prepared the figures and tables and drafted the manuscript. X.W., Y.S., K.W.L.,
680 Z.X.Z., C.S.L., and X.G.W. interpreted the data and discussed the results. X.W. and K.W.L. revised
681 the manuscript. X.W., K.W.L., and C.M.T. provided funding and material support. X.W. and K.W.L.
682 supervised the study. All authors have read and approved the final version of the manuscript.

683 References

- 684 Alpert T, Straube K, Oesterreich FC, Herzl L, Neugebauer KM. 2020. Widespread Transcriptional
685 Readthrough Caused by Nab2 Depletion Leads to Chimeric Transcripts with Retained Introns.
686 *Cell Rep* **33**: 108496.
- 687 Bardia A, Rugo HS, Tolaney SM, Loirat D, Punie K, Oliveira M, Brufsky A, Kalinsky K, Cortes J,
688 Shaughnessy JO et al. 2024. Final Results From the Randomized Phase III ASCENT Clinical
689 Trial in Metastatic Triple-Negative Breast Cancer and Association of Outcomes by Human
690 Epidermal Growth Factor Receptor 2 and Trophoblast Cell Surface Antigen 2 Expression. *J*
691 *Clin Oncol* **42**: 1738-1744.
- 692 Bhattacharya A, Vo DD, Jops C, Kim M, Wen C, Hervoso JL, Pasaniuc B, Gandal MJ. 2023. Isoform-
693 level transcriptome-wide association uncovers genetic risk mechanisms for neuropsychiatric
694 disorders in the human brain. *Nat Genet* **55**: 2117-2128.
- 695 Boixareu C, Taha T, Venkadakrishnan VB, de Bono J, Beltran H. 2025. Targeting the tumour cell
696 surface in advanced prostate cancer. *Nat Rev Urol* doi:10.1038/s41585-025-01014-w.
- 697 Borgers JSW, Lenkala D, Kohler V, Jackson EK, Linssen MD, Hymson S, McCarthy B, O'Reilly
698 Cosgrove E, Balogh KN, Esaulova E et al. 2025. Personalized, autologous neoantigen-specific
699 T cell therapy in metastatic melanoma: a phase 1 trial. *Nat Med* **31**: 881-893.
- 700 Bradley RK, Anczukow O. 2023. RNA splicing dysregulation and the hallmarks of cancer. *Nat Rev*
701 *Cancer* **23**: 135-155.
- 702 Bruce JP, To KF, Lui VWY, Chung GTY, Chan YY, Tsang CM, Yip KY, Ma BBY, Woo JKS, Hui
703 EP et al. 2022. Author Correction: Whole-genome profiling of nasopharyngeal carcinoma
704 reveals viral-host co-operation in inflammatory NF-kappaB activation and immune escape.
705 *Nat Commun* **13**: 4353.

- 706 Burbage M, Rocanin-Arjo A, Baudon B, Arribas YA, Merlotti A, Rookhuizen DC, Heurtebise-
707 Chretien S, Ye M, Houy A, Burgdorf N et al. 2023. Epigenetically controlled tumor antigens
708 derived from splice junctions between exons and transposable elements. *Sci Immunol* **8**:
709 eabm6360.
- 710 Cazes A, Childers BG, Esparza E, Lowy AM. 2022. The MST1R/RON Tyrosine Kinase in Cancer:
711 Oncogenic Functions and Therapeutic Strategies. *Cancers (Basel)* **14**.
- 712 Charlet A, Kappenstein M, Keye P, Klasener K, Endres C, Poggio T, Gorantla SP, Kreutmair S,
713 Sanger J, Illert AL et al. 2022. The IL-3, IL-5, and GM-CSF common receptor beta chain
714 mediates oncogenic activity of FLT3-ITD-positive AML. *Leukemia* **36**: 701-711.
- 715 Chow LK-Y, Chung DL-S, Tao L, Chan KF, Tung SY, Ngan RKC, Ng WT, Lee AW-M, Yau CC,
716 Kwong DL-W. 2022. Epigenomic landscape study reveals molecular subtypes and EBV-
717 associated regulatory epigenome reprogramming in nasopharyngeal carcinoma. *EBioMedicine*
718 **86**.
- 719 Dai W, Zheng H, Cheung AK, Tang CS, Ko JM, Wong BW, Leong MM, Sham PC, Cheung F,
720 Kwong DL et al. 2016. Whole-exome sequencing identifies MST1R as a genetic susceptibility
721 gene in nasopharyngeal carcinoma. *Proc Natl Acad Sci U S A* **113**: 3317-3322.
- 722 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR.
723 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15-21.
- 724 Du W, Liu X, Yang M, Wang W, Sun J. 2021. The Regulatory Role of PRRX1 in Cancer Epithelial-
725 Mesenchymal Transition. *Oncotargets Ther* **14**: 4223-4229.
- 726 Dunn GP, Sherpa N, Manyanga J, Johanns TM. 2022. Considerations for personalized neoantigen
727 vaccination in Malignant glioma. *Adv Drug Deliv Rev* **186**: 114312.
- 728 Gallego-Paez LM, Edwards WJS, Chanduri M, Guo Y, Koorman T, Lee CY, Grexa N, Derksen P,
729 Yan J, Schwartz MA et al. 2023. TLN1 contains a cancer-associated cassette exon that alters
730 talin-1 mechanosensitivity. *J Cell Biol* **222**.
- 731 Glinos DA, Garborcauskas G, Hoffman P, Ehsan N, Jiang L, Gokden A, Dai X, Aguet F, Brown KL,
732 Garimella K et al. 2022. Transcriptome variation in human tissues revealed by long-read
733 sequencing. *Nature* **608**: 353-359.

- 734 Hallgren J, Tsirigos KD, Pedersen MD, Almagro Armenteros JJ, Marcatili P, Nielsen H, Krogh A,
735 Winther O. 2022. DeepTMHMM predicts alpha and beta transmembrane proteins using deep
736 neural networks. *bioRxiv* doi:10.1101/2022.04.08.487609: 2022.2004.2008.487609.
- 737 Hong M, Tao S, Zhang L, Diao LT, Huang X, Huang S, Xie SJ, Xiao ZD, Zhang H. 2020. RNA
738 sequencing: new technologies and applications in cancer research. *J Hematol Oncol* **13**: 166.
- 739 Hook PW, Timp W. 2023. Beyond assembly: the increasing flexibility of single-molecule sequencing
740 technology. *Nat Rev Genet* **24**: 627-641.
- 741 Hu W, Wu Y, Shi Q, Wu J, Kong D, Wu X, He X, Liu T, Li S. 2022. Systematic characterization of
742 cancer transcriptome at transcript resolution. *Nat Commun* **13**: 6803.
- 743 Hu Z, Guo X, Li Z, Meng Z, Huang S. 2024. The neoantigens derived from transposable elements - A
744 hidden treasure for cancer immunotherapy. *Biochim Biophys Acta Rev Cancer* **1879**: 189126.
- 745 Hu Z, Yuan J, Long M, Jiang J, Zhang Y, Zhang T, Xu M, Fan Y, Tanyi JL, Montone KT et al. 2021.
746 The Cancer Surfaceome Atlas integrates genomic, functional and drug response data to
747 identify actionable targets. *Nat Cancer* **2**: 1406-1422.
- 748 Huang KK, Huang J, Wu JKL, Lee M, Tay ST, Kumar V, Ramnarayanan K, Padmanabhan N, Xu C,
749 Tan ALK et al. 2021. Long-read transcriptome sequencing reveals abundant promoter
750 diversity in distinct molecular subtypes of gastric cancer. *Genome Biol* **22**: 44.
- 751 Joglekar A, Hu W, Zhang B, Narykov O, Diekhans M, Marrocco J, Balacco J, Ndhlovu LC, Milner
752 TA, Fedrigo O et al. 2024. Single-cell long-read sequencing-based mapping reveals
753 specialized splicing patterns in developing and adult mouse and human brain. *Nat Neurosci* **27**:
754 1051-1063.
- 755 Jurtz V, Paul S, Andreatta M, Marcatili P, Peters B, Nielsen M. 2017. NetMHCpan-4.0: Improved
756 Peptide-MHC Class I Interaction Predictions Integrating Eluted Ligand and Peptide Binding
757 Affinity Data. *J Immunol* **199**: 3360-3368.
- 758 Kahles A, Lehmann KV, Toussaint NC, Huser M, Stark SG, Sachsenberg T, Stegle O, Kohlbacher O,
759 Sander C, Cancer Genome Atlas Research N et al. 2018. Comprehensive Analysis of
760 Alternative Splicing Across Tumors from 8,705 Patients. *Cancer Cell* **34**: 211-224 e216.

- 761 Kam NW, Lau CY, Lau JYH, Dai X, Liang Y, Lai SPH, Chung MKY, Yu VZ, Qiu W, Yang M et al.
762 2025. Cell-associated galectin 9 interacts with cytotoxic T cells confers resistance to tumor
763 killing in nasopharyngeal carcinoma through autophagy activation. *Cell Mol Immunol* **22**:
764 260-281.
- 765 Kechin A, Boyarskikh U, Kel A, Filipenko M. 2017. cutPrimers: A New Tool for Accurate Cutting of
766 Primers from Reads of Targeted Next Generation Sequencing. *J Comput Biol* **24**: 1138-1143.
- 767 Kim S, Pevzner PA. 2014. MS-GF+ makes progress towards a universal database search tool for
768 proteomics. *Nat Commun* **5**: 5277.
- 769 Klibi J, Niki T, Riedel A, Pioche-Durieu C, Souquere S, Rubinstein E, Le Moulec S, Guigay J,
770 Hirashima M, Guemira F et al. 2009. Blood diffusion and Th1-suppressive effects of galectin-
771 9-containing exosomes released by Epstein-Barr virus-infected nasopharyngeal carcinoma
772 cells. *Blood* **113**: 1957-1966.
- 773 Kumar H, Luo R, Wen J, Yang C, Zhou X, Kim P. 2024. FusionNeoAntigen: a resource of fusion
774 gene-specific neoantigens. *Nucleic Acids Res* **52**: D1276-D1288.
- 775 Lauss M, Donia M, Harbst K, Andersen R, Mitra S, Rosengren F, Salim M, Vallon-Christersson J,
776 Torngren T, Kvist A et al. 2020. Author Correction: Mutational and putative neoantigen load
777 predict clinical benefit of adoptive T cell therapy in melanoma. *Nat Commun* **11**: 1714.
- 778 Li B, Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-Seq data with or without
779 a reference genome. *BMC Bioinformatics* **12**: 323.
- 780 Liberzon A, Birger C, Thorvaldsdottir H, Ghandi M, Mesirov JP, Tamayo P. 2015. The Molecular
781 Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* **1**: 417-425.
- 782 Lo AK, Dawson CW, Lung HL, Wong KL, Young LS. 2021. The Role of EBV-Encoded LMP1 in the
783 NPC Tumor Microenvironment: From Function to Therapy. *Front Oncol* **11**: 640207.
- 784 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-seq
785 data with DESeq2. *Genome Biol* **15**: 550.
- 786 Lu SX, De Neef E, Thomas JD, Sabio E, Rousseau B, Gigoux M, Knorr DA, Greenbaum B, Elhanati
787 Y, Hogg SJ et al. 2021. Pharmacologic modulation of RNA splicing enhances anti-tumor
788 immunity. *Cell* **184**: 4032-4047 e4031.

- 789 Maurin M, Ranjouri M, Megino-Luque C, Newberg JY, Du D, Martin K, Miner RE, 3rd, Prater MS,
790 Wee DKB, Centeno B et al. 2023. RBFOX2 deregulation promotes pancreatic cancer
791 progression and metastasis through alternative splicing. *Nat Commun* **14**: 8444.
- 792 Merlotti A, Sadacca B, Arribas YA, Ngoma M, Burbage M, Goudot C, Houy A, Rocanin-Arjo A,
793 Lalanne A, Seguin-Givelet A et al. 2023. Noncanonical splicing junctions between exons and
794 transposable elements represent a source of immunogenic recurrent neo-antigens in patients
795 with lung cancer. *Sci Immunol* **8**: eabm6359.
- 796 Monzo C, Liu T, Conesa A. 2025. Transcriptomics in the era of long-read sequencing. *Nat Rev Genet*
797 doi:10.1038/s41576-025-00828-z.
- 798 Moon Y, Herrmann CJ, Mironov A, Zavolan M. 2025. PolyASite v3.0: a multi-species atlas of
799 polyadenylation sites inferred from single-cell RNA-sequencing data. *Nucleic Acids Res* **53**:
800 D197-D204.
- 801 Naro C, Antonioni A, Medici V, Caggiano C, Jolly A, de la Grange P, Bielli P, Paronetto MP, Sette C.
802 2024. Splicing targeting drugs highlight intron retention as an actionable vulnerability in
803 advanced prostate cancer. *J Exp Clin Cancer Res* **43**: 58.
- 804 Neri P, Leblay N, Lee H, Gulla A, Bahlis NJ, Anderson KC. 2024. Just scratching the surface: novel
805 treatment approaches for multiple myeloma targeting cell membrane proteins. *Nat Rev Clin*
806 *Oncol* **21**: 590-609.
- 807 Noguchi S, Arakawa T, Fukuda S, Furuno M, Hasegawa A, Hori F, Ishikawa-Kato S, Kaida K, Kaiho
808 A, Kanamori-Katayama M. 2017. FANTOM5 CAGE profiles of human and mouse samples.
809 *Scientific data* **4**: 170112.
- 810 Odum MT, Teufel F, Thumuluri V, Almagro Armenteros JJ, Johansen AR, Winther O, Nielsen H.
811 2024. DeepLoc 2.1: multi-label membrane protein type prediction using protein language
812 models. *Nucleic Acids Res* **52**: W215-W220.
- 813 Ott PA, Hu-Lieskovan S, Chmielowski B, Govindan R, Naing A, Bhardwaj N, Margolin K, Awad
814 MM, Hellmann MD, Lin JJ et al. 2020. A Phase Ib Trial of Personalized Neoantigen Therapy
815 Plus Anti-PD-1 in Patients with Advanced Melanoma, Non-small Cell Lung Cancer, or
816 Bladder Cancer. *Cell* **183**: 347-362 e324.

- 817 Pan Y, Phillips JW, Zhang BD, Noguchi M, Kutschera E, McLaughlin J, Nesterenko PA, Mao Z,
818 Bangayan NJ, Wang R et al. 2023. IRIS: Discovery of cancer immunotherapy targets arising
819 from pre-mRNA alternative splicing. *Proc Natl Acad Sci U S A* **120**: e2221116120.
- 820 Pardo-Palacios FJ, Arzalluz-Luque A, Kondratova L, Salguero P, Mestre-Tomas J, Amorin R,
821 Estevan-Morio E, Liu T, Nanni A, McIntyre L et al. 2024. SQANTI3: curation of long-read
822 transcriptomes for accurate identification of known and novel isoforms. *Nat Methods* **21**: 793-
823 797.
- 824 Pertea G, Pertea M. 2020. GFF Utilities: GffRead and GffCompare. *F1000Res* **9**.
- 825 Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. 2015. StringTie enables
826 improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol* **33**: 290-295.
- 827 Platten M, Bunse L, Wick A, Bunse T, Le Cornet L, Harting I, Sahm F, Sanghvi K, Tan CL, Poschke I
828 et al. 2021. A vaccine targeting mutant IDH1 in newly diagnosed glioma. *Nature* **592**: 463-
829 468.
- 830 Sethna Z, Guasp P, Reiche C, Milighetti M, Ceglia N, Patterson E, Lihm J, Payne G, Lyudovyk O,
831 Rojas LA et al. 2025. RNA neoantigen vaccines prime long-lived CD8(+) T cells in pancreatic
832 cancer. *Nature* **639**: 1042-1051.
- 833 Siak PY, Heng WS, Teoh SSH, Lwin YY, Cheah SC. 2023. Precision medicine in nasopharyngeal
834 carcinoma: comprehensive review of past, present, and future prospect. *J Transl Med* **21**: 786.
- 835 Smart AC, Margolis CA, Pimentel H, He MX, Miao D, Adeegbe D, Fugmann T, Wong KK, Van
836 Allen EM. 2018. Intron retention is a source of neoepitopes in cancer. *Nat Biotechnol* **36**:
837 1056-1058.
- 838 Somalraju S, Salem DH, Janga SC. 2025. Investigating RNA dynamics from single molecule
839 transcriptomes. *Trends Genet* doi:10.1016/j.tig.2025.05.001.
- 840 Supek F, Lehner B, Lindeboom RGH. 2021. To NMD or Not To NMD: Nonsense-Mediated mRNA
841 Decay in Cancer and Other Genetic Diseases. *Trends Genet* **37**: 657-668.
- 842 Szolek A, Schubert B, Mohr C, Sturm M, Feldhahn M, Kohlbacher O. 2014. OptiType: precision HLA
843 typing from next-generation sequencing data. *Bioinformatics* **30**: 3310-3316.

- 844 Trincado JL, Entizne JC, Hysenaj G, Singh B, Skalic M, Elliott DJ, Eyraas E. 2018. SUPPA2: fast,
845 accurate, and uncertainty-aware differential splicing analysis across multiple conditions.
846 *Genome Biol* **19**: 40.
- 847 Tsang CM, Lui VWY, Bruce JP, Pugh TJ, Lo KW. 2020. Translational genomics of nasopharyngeal
848 cancer. *Semin Cancer Biol* **61**: 84-100.
- 849 Wang X, Yu L, Zhou X, Chung GT, Liu AM, Chan YY, Wu M, Chau KY, Lo KW, Wu AR. 2025.
850 Characterizing resistant cellular states in nasopharyngeal carcinoma during EBV lytic
851 induction. *Oncogene* **44**: 1805-1819.
- 852 Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, Feng T, Zhou L, Tang W, Zhan L et al. 2021.
853 clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. *Innovation (Camb)*
854 **2**: 100141.
- 855 Xie H, Zhang K, Yin H, Zhang S, Pan S, Wu R, Han Y, Xu Y, Jiang W, You B. 2025.
856 Acetyltransferase NAT10 inhibits T-cell immunity and promotes nasopharyngeal carcinoma
857 progression through DDX5/HMGB1 axis. *J Immunother Cancer* **13**.
- 858 Xie N, Shen G, Gao W, Huang Z, Huang C, Fu L. 2023. Neoantigens: promising targets for cancer
859 therapy. *Signal Transduct Target Ther* **8**: 9.
- 860 Zaidi S, Park J, Chan JM, Roudier MP, Zhao JL, Gopalan A, Wadosky KM, Patel RA, Sayar E,
861 Karthaus WR et al. 2024. Single-cell analysis of treatment-resistant prostate cancer:
862 Implications of cell state changes for cell surface antigen-targeted therapies. *Proc Natl Acad*
863 *Sci U S A* **121**: e2322203121.
- 864 Zelba H, Shao B, Rabsteyn A, Reinhardt A, Greve C, Oenning L, Kayser S, Kyzirakos C, Latzer P,
865 Riedlinger T et al. 2025. In-depth characterization of vaccine-induced neoantigen-specific T
866 cells in patients with IDH1-mutant glioma undergoing personalized peptide vaccination. *J*
867 *Immunother Cancer* **13**.
- 868

869 **Figure legends**870 **Fig. 1. Profiling of full-length transcripts in nasopharyngeal carcinoma by long-read sequencing.**

871 (A) Pie chart illustrating the distribution of full-length transcripts across different subcategories. A
 872 total of 51,115 full-length transcripts were identified across 18 long-read sequencing samples, with the
 873 numbers and proportions of different transcript subcategories highlighted in different colors. (B) Bar
 874 plots showing LR-seq isoforms detected in individual NPC tumor samples and immortalized NP cell
 875 lines, colored by the subcategories defined in (A) and grouped by sample origin. (C) Histograms
 876 comparing length distributions of long-read transcripts between different subcategories. (D) Bar plot
 877 showing the number of transcripts identified per gene by LR-seq. (E) Boxplot comparing the
 878 proportions of long-read transcripts supported by short-read Illumina RNA-seq among different
 879 subcategories in matched samples.

880 **Fig. 2. Characterization of full-length transcripts identified in LR-seq.** (A) Barplot comparing the

881 coding potential, predicted NMD status, and splice site canonicity between transcript subcategories.
 882 (B) Boxplots comparing isoform length, CDS length, and junction coverage across FSM, NIC and
 883 NNC transcript categories. (C) Dot plot showing pathways significantly enriched among genes with
 884 novel transcripts detected by LR-seq ($FDR < 0.05$), based on the MSigDB C2, C5, and Hallmark
 885 collections. Enrichment remained significant after controlling for gene expression using expression-
 886 matched permutation analysis. (D) Dot plot illustrating the numbers of novel LR-seq transcripts
 887 compared with annotated GENCODE isoforms detected in selected oncogenes (left) and tumor
 888 suppressors (right). (E) Pie chart showing the proportion of novel transcripts sharing the same
 889 translation start site as annotated protein isoforms. NS, not significant; $*P < 0.05$, $**P < 0.01$, $***P <$
 890 0.001 , $****P < 0.0001$.

891 **Fig. 3. Novel transcripts contribute to isoform-level diversity and differential alternative splicing**

892 **events.** (A) Schematic illustration of the seven types of AS events, including A5: Alternative 5' splice
 893 site, A3: Alternative 3' splice site, AF: Alternative first exon, AL: Alternative last exon, MX:
 894 Mutually exclusive exons, RI: Retained intron, and SE: Skipping exon. (B) Bar plot comparing
 895 alternative splicing events before and after incorporating novel transcripts into GENCODE v38

896 transcriptome. (C) Dot plot showing differential alternative splicing events between NPC tumor
897 samples and immortalized NP cell lines, identified by SUPPA2 with $P < 0.05$ and $|\Delta\text{PSI}| > 0.1$. (D)
898 Schematic illustrations of an example of a differential alternative splicing event involving the
899 inclusion of a cassette exon (17b) in *TLNI* and the resulting transcripts. (E) Bar plot comparing
900 differential isoform usage of *TLNI* (left), with and without exon 17b, between NPC samples and
901 immortalized NP cell lines, validated by RT-qPCR (right). Each RT-qPCR experiment was repeated
902 three times with independent RT preparations, with the standard deviation of each triplicate indicated
903 by the error bars.

904 **Fig. 4. Identification and validation of tumor-specific transcripts.** (A) Density curves comparing
905 transcript expression across different transcript subcategories. (B) Bar plot showing the numbers of
906 FSM, NIC and NNC transcripts across differential median expression thresholds. (C) Volcano plot
907 illustrating differential transcript expression analysis between NPC tumor samples and immortalized
908 NP cell lines (BH-adjusted $P < 0.05$ and $|\log_2\text{FC}| > 2$). (D) Ring charts showing the classification of
909 transcripts according to their expression specificity. (E) Heatmap showing the expression of a subset
910 of predicted tumor-specific novel transcripts identified in LR-seq. (F) Bar plots depicting RT-qPCR
911 validation results of transcripts in (E). RT-qPCR assays were carried out in three independent RT
912 replicates, with error bars representing the standard deviation.

913 **Fig. 5. Identification of potential neojunction-derived neoantigens from tumor-specific novel**
914 **transcripts using TS-SNAD.** (A) Pie chart showing the distribution of splice junction categories
915 detected by LR-seq, with known splice junctions (blue) and novel splice junctions (yellow). (B) Bar
916 plot showing the number of transcripts with different numbers of novel splice junctions. Inset pie chart
917 depicting the proportion of canonical and non-canonical splice sites. (C) Bar plot showing the
918 distribution of transcripts containing novel splice junctions, categorized as NNC, NIC, intergenic,
919 genic, fusion, or antisense. (D) Schematic overview of the TS-SNAD computational workflow for
920 identifying neojunction-derived neoantigens from tumor-specific novel transcripts. (E) Ring plot
921 showing that a total of 3,326 neojunction-derived neoantigens were identified, of which 533 were
922 predicted to bind strongly to at least one NPC HLA class I allele. (F) Bar plot showing the distribution
923 of 533 strong-binding neoantigens by the number of tumor-specific novel transcripts from which they

924 were derived. (G) Bar plot showing the number of distinct HLA class I alleles bound by the 533
925 neoantigens.

926 **Fig. 6. Neojunction-derived neoantigens effectively elicit immunogenic responses.** (A) ELISpot
927 assay results demonstrating that six out of the 34 splicing-derived neoantigens elicited immunogenic
928 responses. PBMCs from three independent HLA-B*40:01 donors were stimulated with peptides (20
929 $\mu\text{g}/\text{mL}$), and IFN- γ secretion was measured by ELISpot ($n = 3$). DMSO served as the negative control,
930 and anti-CD3 mAb as the positive control. Data are presented as mean \pm SD. (B) Schematic
931 illustrations of the exon-intron structures of the novel transcripts from which the six immunogenic
932 neoantigens were derived.

933 **Fig. 7. Proteomic and protein expression validation of ORFs encoding immunogenic**
934 **neojunction-derived neoantigens in NPC.** (A) Mass spectrometry provided protein-level evidence
935 supporting the existence of neojunction-derived neoantigens. (B) Western blot analysis showing the
936 expression of FLAG-tagged ORFs encoding the validated immunogenic neojunction-derived
937 neoantigens in HEK293T and NP69 cells after transfection.

938

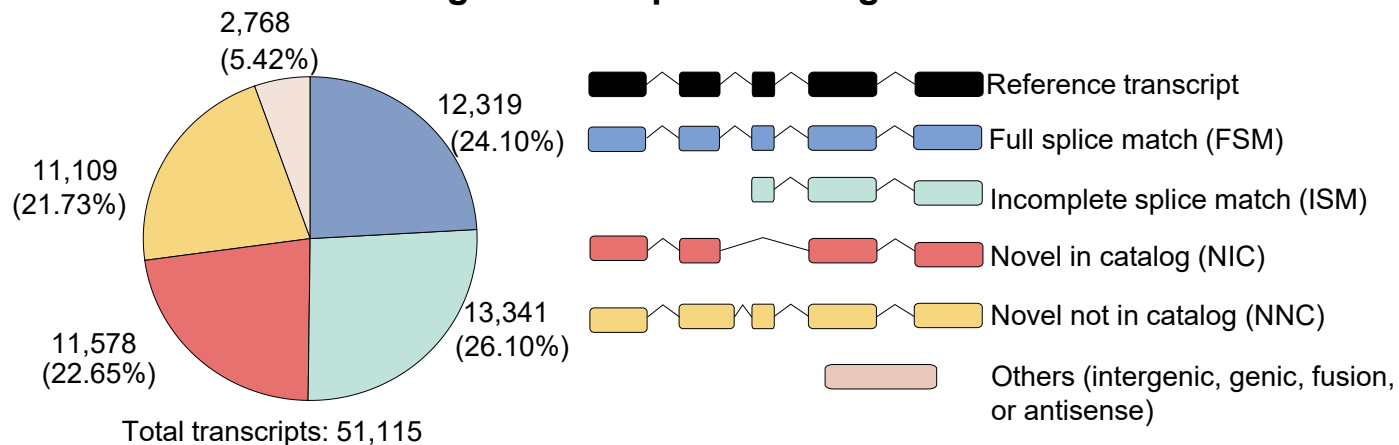
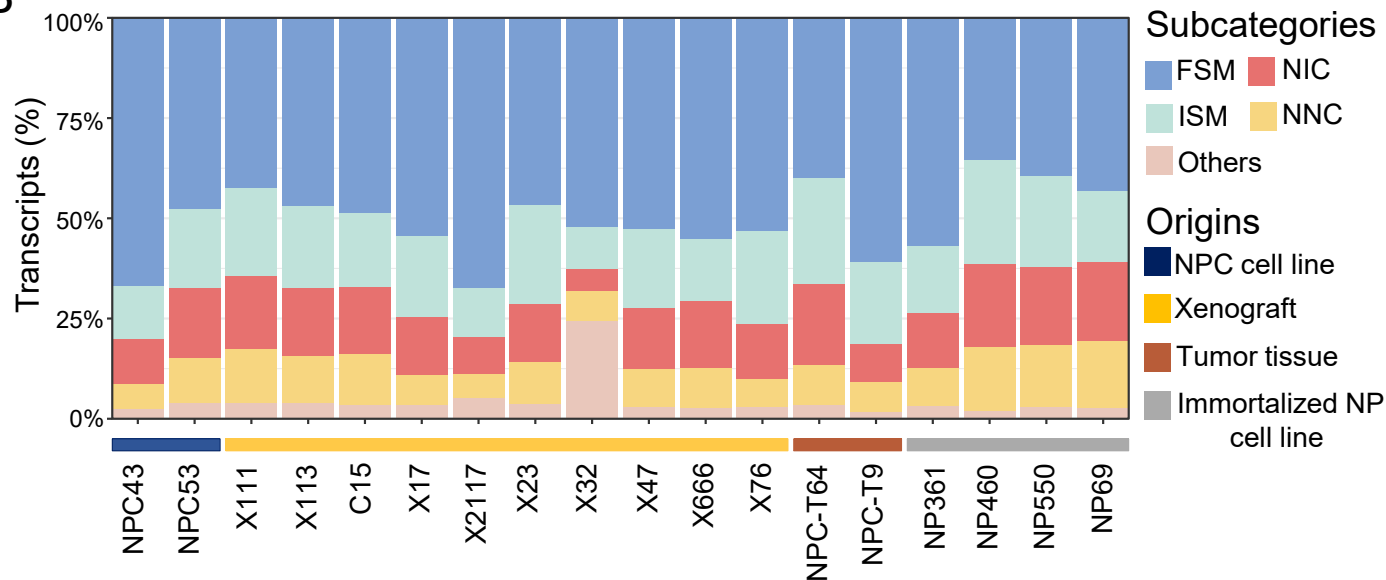
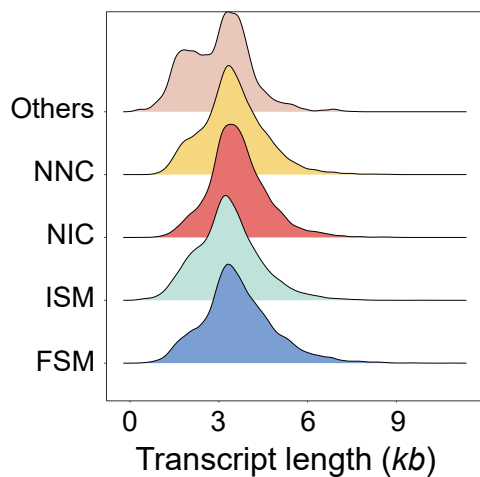
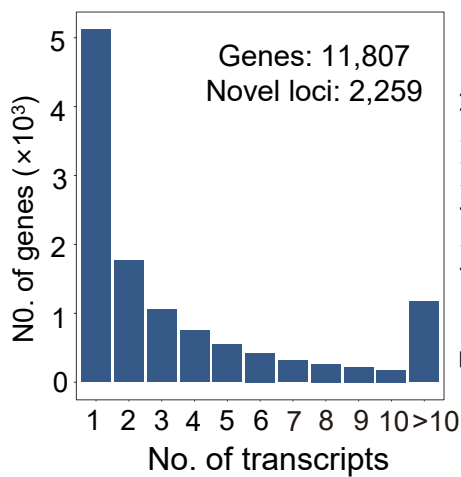
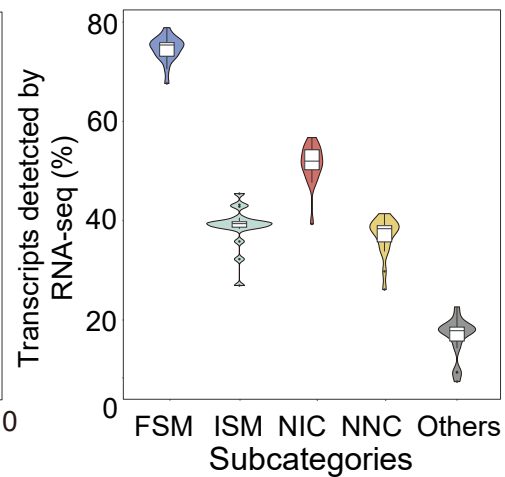
939 Tables

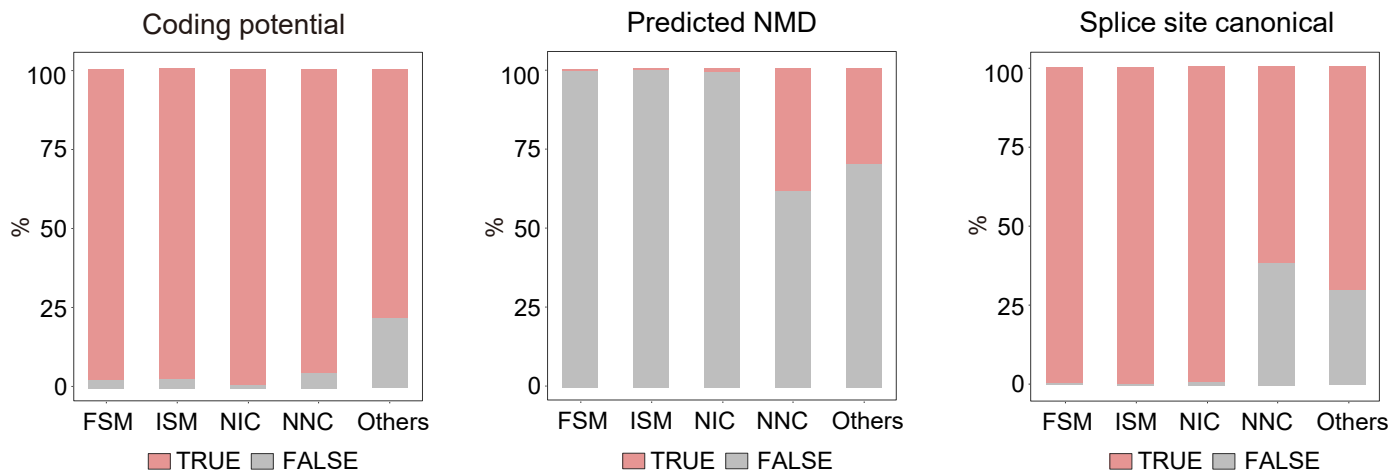
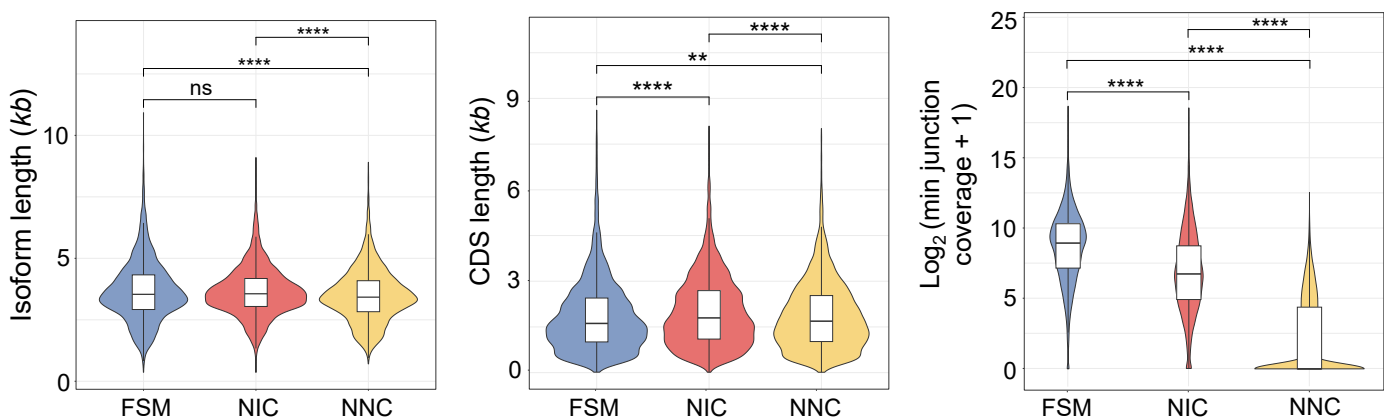
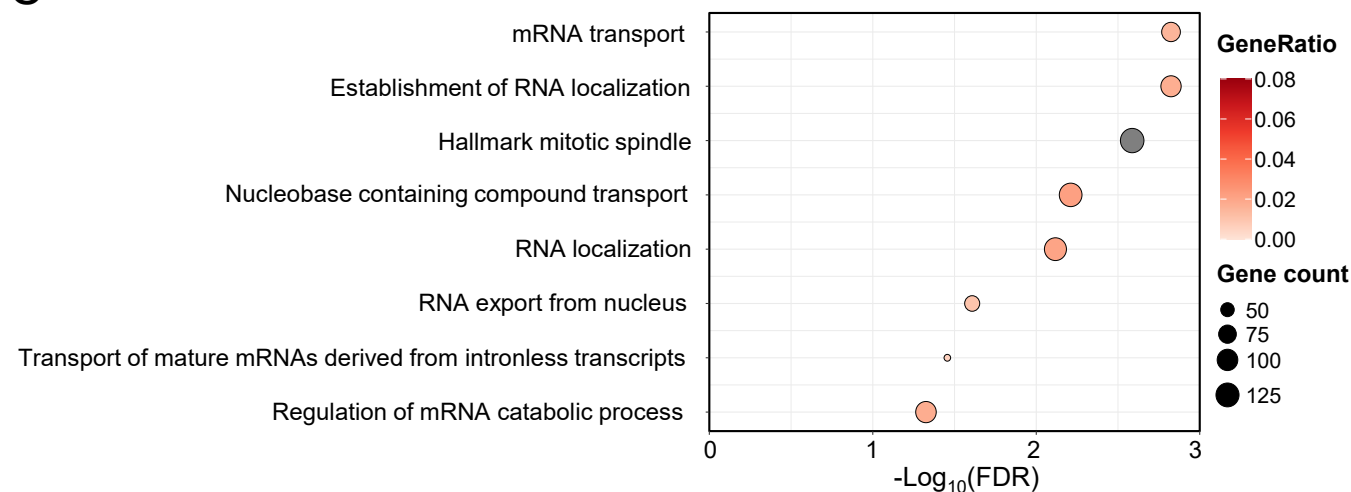
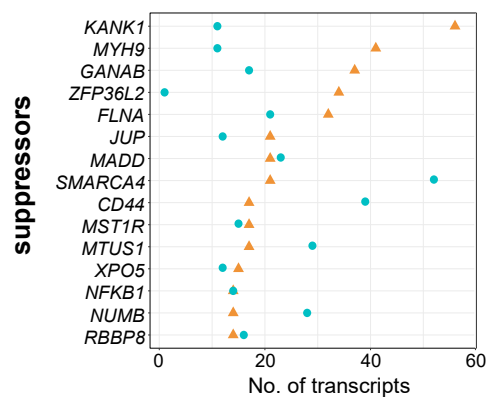
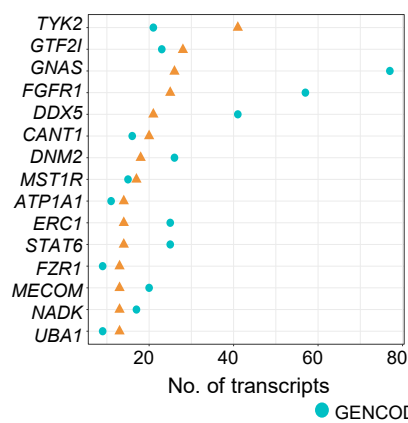
940

Table 1. NPC/immortalized NP cell lines for long-read sequencing

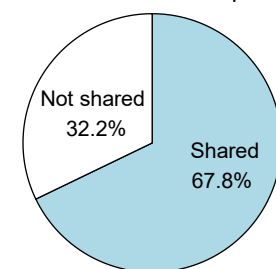
Sample Name	Sample Type	EBV Infection
NPC43	NPC cell line	Yes
NPC53	NPC cell line	No
NP361	Immortalized nasopharyngeal epithelial (NP) cell lines	No
NP460	Immortalized nasopharyngeal epithelial (NP) cell lines	No
NP550	Immortalized nasopharyngeal epithelial (NP) cell lines	No
NP69	Immortalized nasopharyngeal epithelial (NP) cell lines	No
X111	Patient-derived tumor xenograft model	Yes
X113	Patient-derived tumor xenograft model	Yes
C15	Patient-derived tumor xenograft model	Yes
X17	Patient-derived tumor xenograft model	Yes
X2117	Patient-derived tumor xenograft model	Yes
X23	Patient-derived tumor xenograft model	Yes
X32	Patient-derived tumor xenograft model	Yes
X47	Patient-derived tumor xenograft model	Yes
X666	Patient-derived tumor xenograft model	Yes
X76	Patient-derived tumor xenograft model	Yes
NPC-T64	Tumor tissue	Yes
NPC-T9	Tumor tissue	Yes

941

A**Full-length transcripts subcategories****B****C****D****E**

A**B****C****D****E**

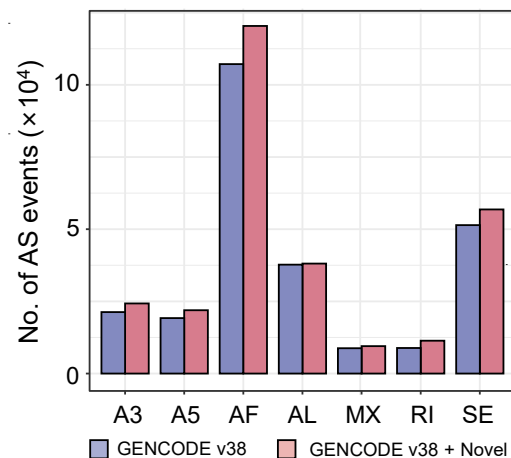
Novel transcripts sharing start codon with known proteins



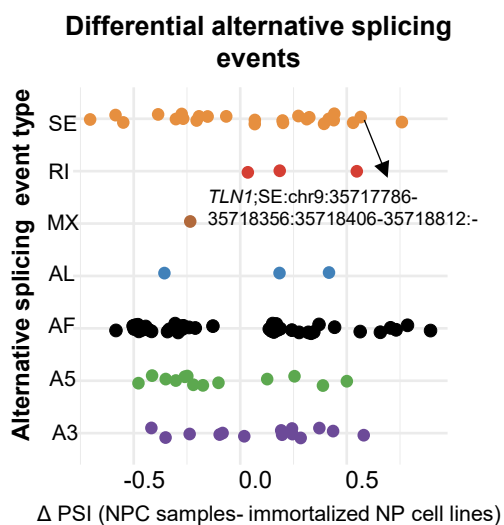
A



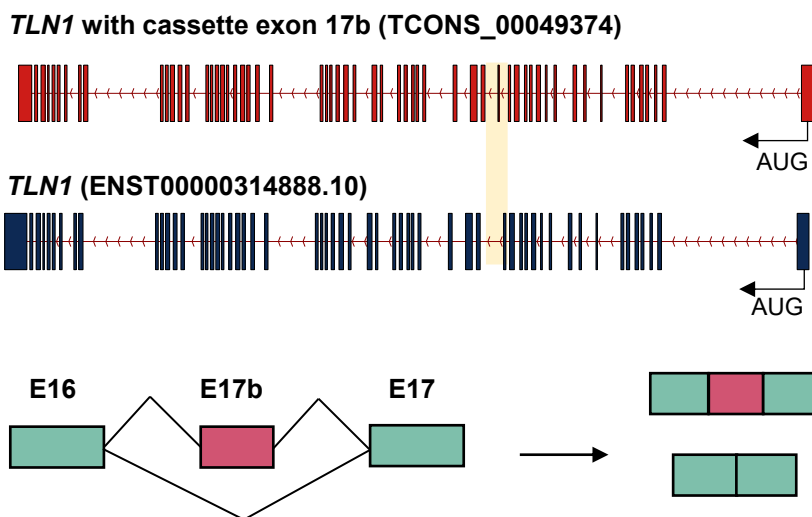
B



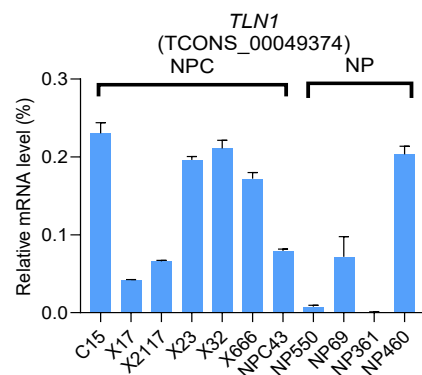
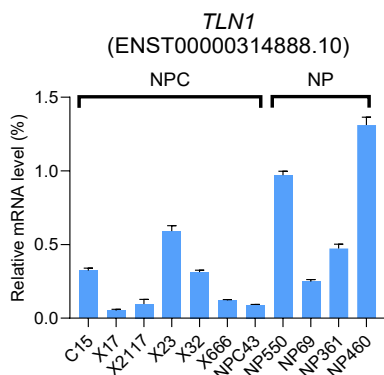
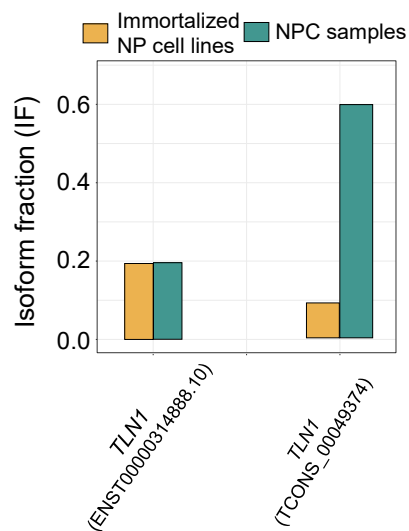
C

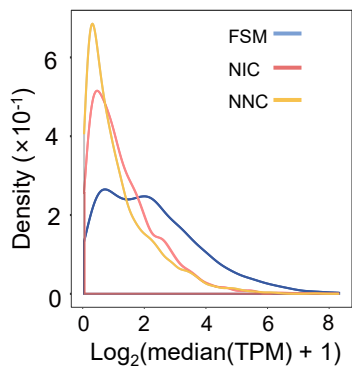
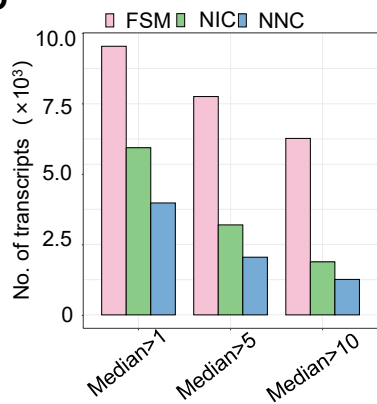
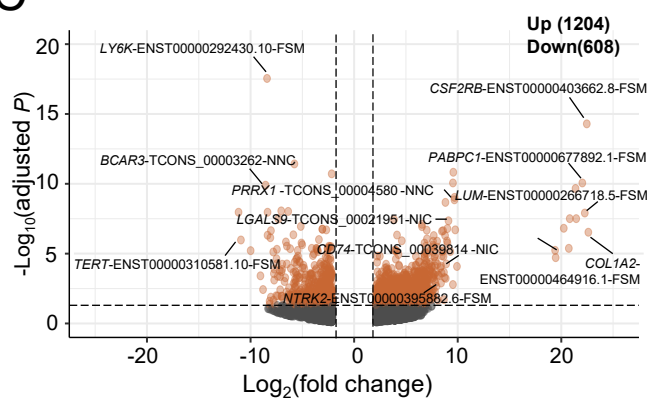
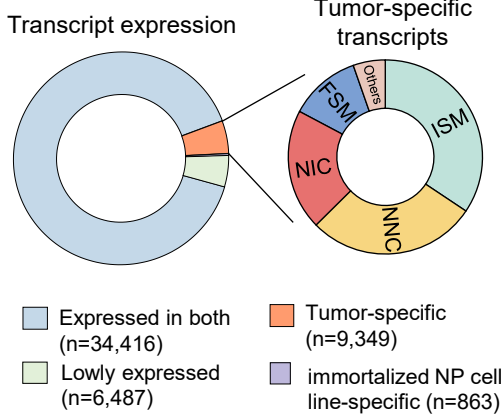
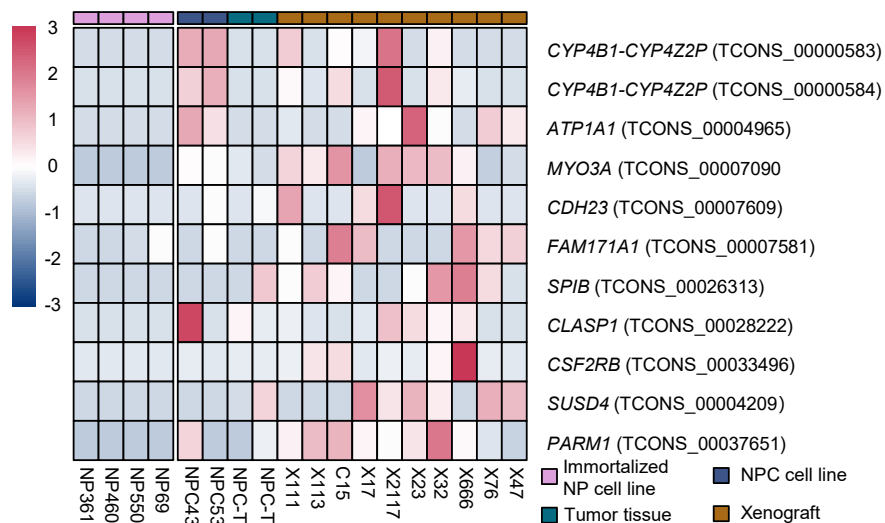
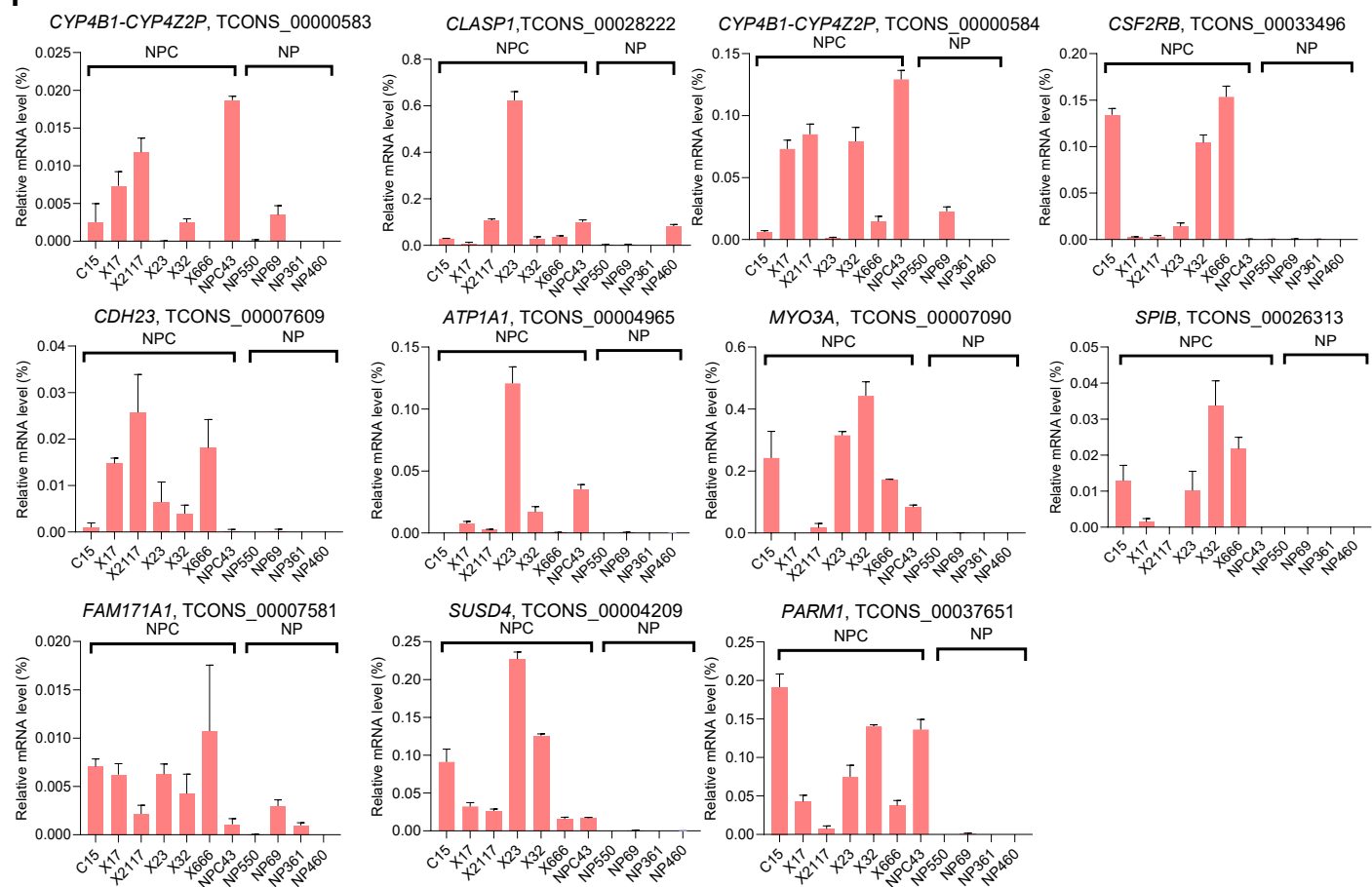


D



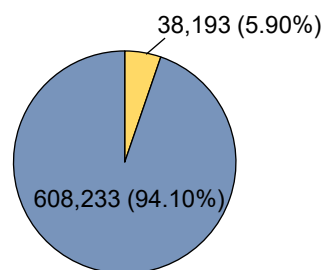
E



A**B****C****D****E****F**

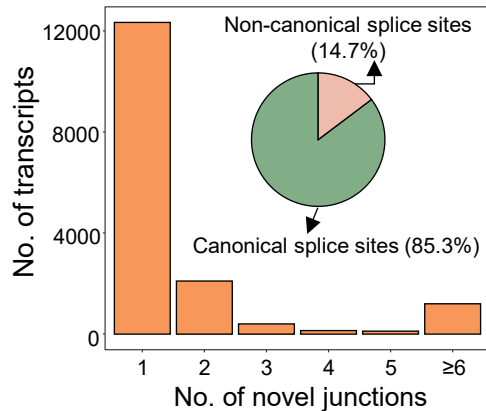
A

Junction categories



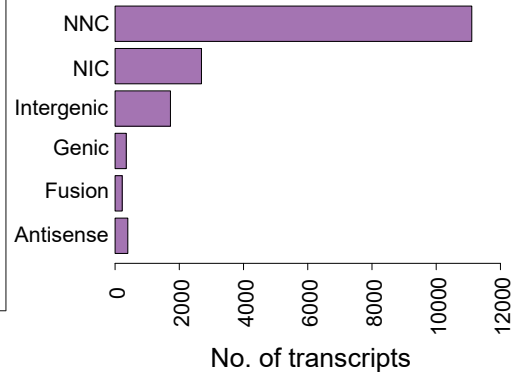
■ Novel splice junctions
■ Known splice junctions

B



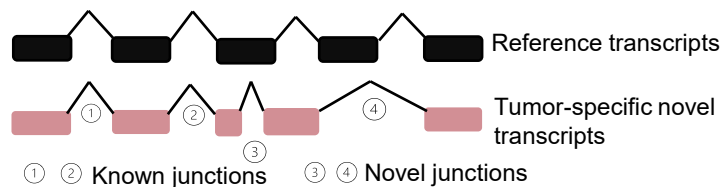
C

Distribution of novel junctions

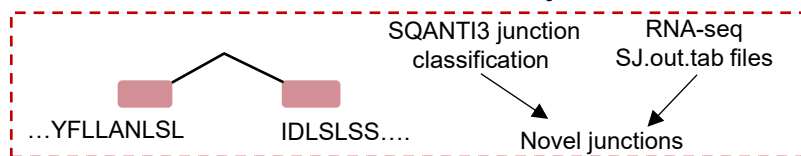


D

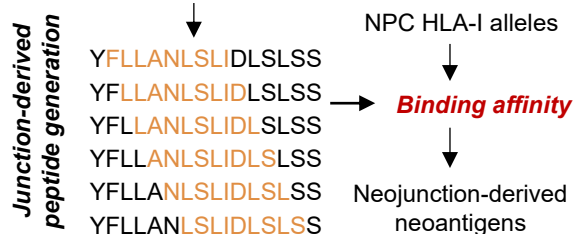
TS-SNAD workflow



Novel junction identification

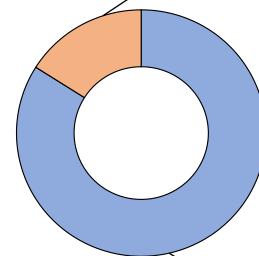


Exclusion of GENCODE v38 annotated or GTEx splice junctions



E

With high HLA-I binding affinity (533, 16%)

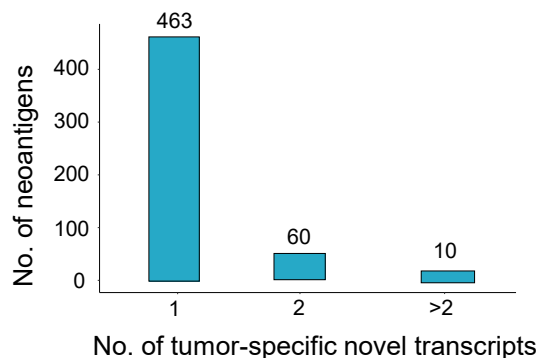


IC ₅₀ binding affinity (nM)	Count
0-20	123
20-50	128
50-100	108
100-200	80
200-300	83
300-500	11

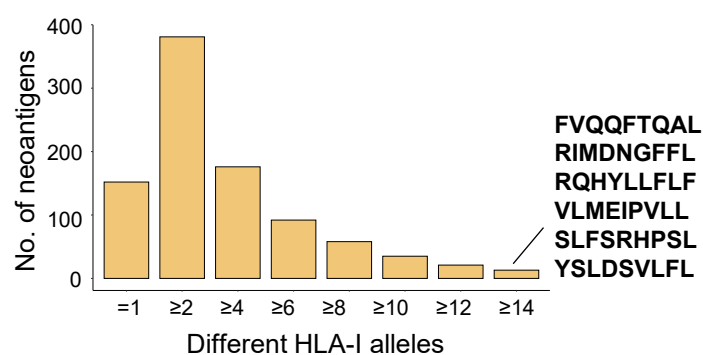
With low HLA-I binding affinity (2,793, 84%)

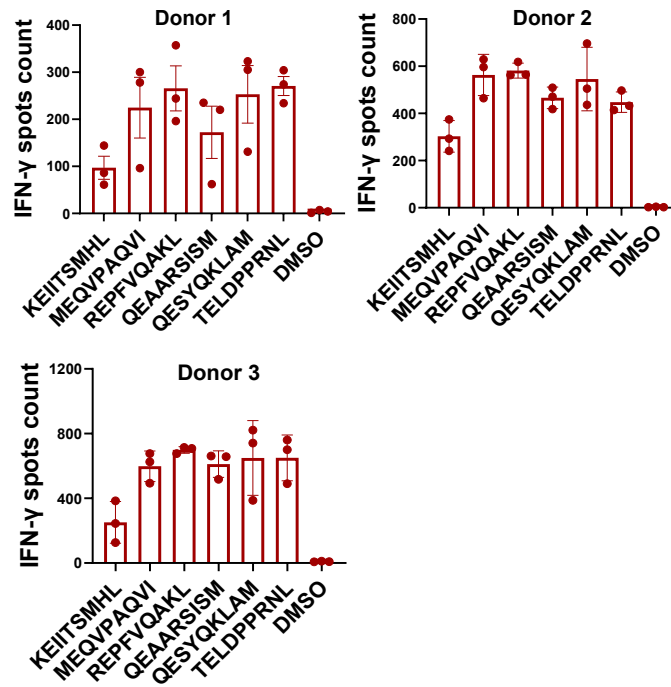
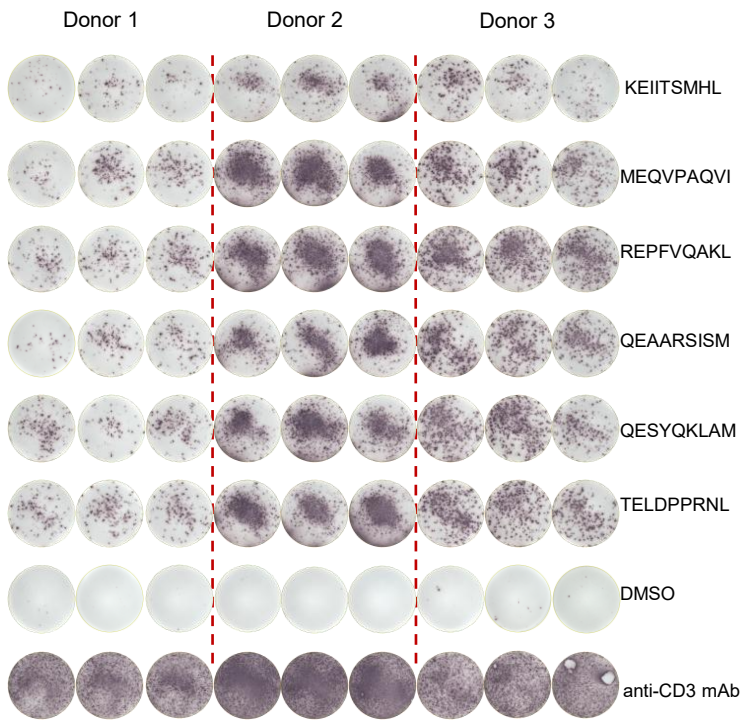
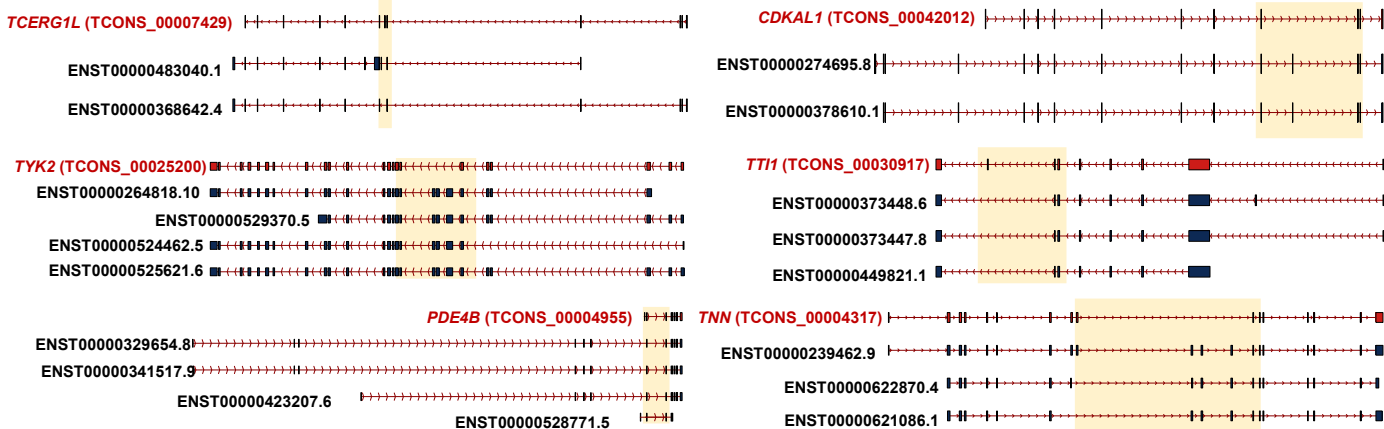
Total: 3,326 neojunction-derived neoantigens

F

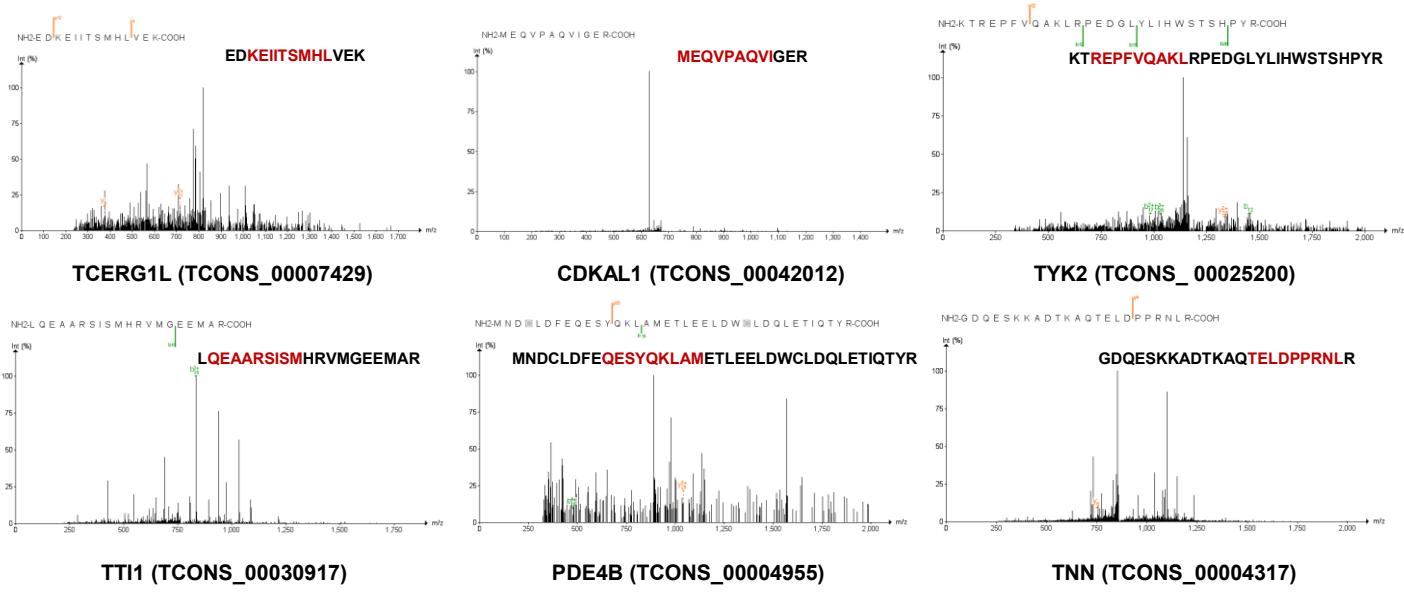


G



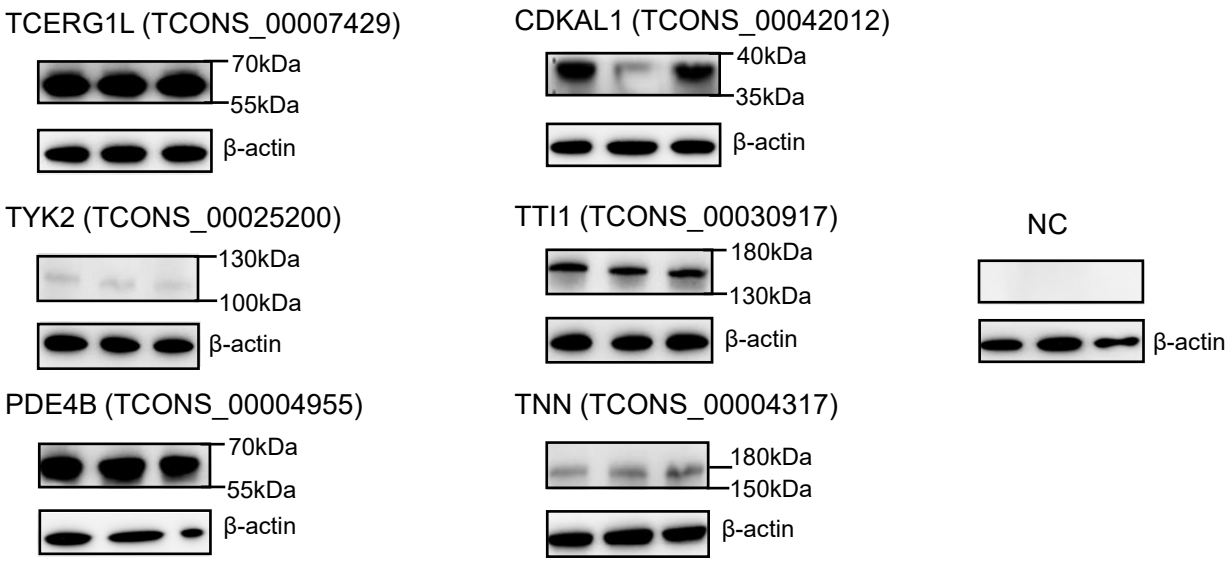
A**B**

A



B

HEK293T



NP69

