



## Dynamics of intronic polyadenylation in the hematopoietic lineage and its regulation by DNA methylation

Richa Rashmi, Abhinaya Muruganandham, Pranita Borkar, et al.

*Genome Res.* published online April 13, 2026

Access the most recent version at doi:[10.1101/gr.281044.125](https://doi.org/10.1101/gr.281044.125)

---

|                                 |  |
|---------------------------------|--|
| <b>P&lt;P</b>                   | Published online April 13, 2026 in advance of the print journal.   |
| <b>Accepted Manuscript</b>      | Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.   |
| <b>Open Access</b>              | Freely available online through the <i>Genome Research</i> Open Access option.   |
| <b>Creative Commons License</b> | 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 <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> . |
| <b>Email Alerting Service</b>   | Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .  |

---



---

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

---

Published by Cold Spring Harbor Laboratory Press

# 1 **Dynamics of intronic polyadenylation in the hematopoietic lineage and its** 2 **regulation by DNA methylation**

3 Richa Rashmi<sup>1</sup>, Abhinaya Muruganandham<sup>1,2</sup>, Pranita Borkar<sup>1</sup>, Sumana Mallick<sup>1,2</sup>, Taylor  
4 Hubbs<sup>1</sup>, Ari Aviles<sup>1,3</sup>, Daniel Chung<sup>1</sup>, Irtisha Singh<sup>1,2,3\*</sup>

5 <sup>1</sup>Department of Cell Biology & Genetics, Texas A&M University Health Science Center, Bryan,  
6 TX, 77807, USA; <sup>2</sup>Department of Biomedical Engineering, Texas A&M University, College  
7 Station, TX, 77843, USA; <sup>3</sup>Interdisciplinary Program in Genetics and Genomics, Texas A&M  
8 University, College Station, TX, 77843, USA

## 9 Abstract

10 Intronic polyadenylation (IPA) is a key mechanism driving transcriptome diversity, yet its detection  
11 and functional characterization remain challenging due to complex splicing patterns and  
12 complexity of intronic regions. Here, we introduce IPAseek, a dynamic programming-based  
13 computational framework that leverages the Pruned Exact Linear Time (PELT) algorithm and  
14 Changepoints Over a Range of Penalties (CROPS) to enable *de novo* identification of IPA events  
15 from bulk RNA-seq data. IPAseek robustly detects both composite and skipped IPA isoforms.

16

17 Applying IPAseek to bulk RNA-seq of hematopoietic cell-types, reveals lineage and stage-specific  
18 IPA signatures, with lymphoid cells exhibiting higher IPA site usage than myeloid cells. Temporal  
19 profiling during megakaryocyte differentiation uncovers dynamic, gene-specific IPA regulation  
20 linked to functional pathways including peroxisomal metabolism and autophagy which are known  
21 to play a crucial role in megakaryocytic differentiation, impacting the development and maturation  
22 of megakaryocytes. Further, integrative analysis demonstrates that IPA site usage is associated  
23 with lower DNA methylation within introns, supporting a regulatory axis connecting epigenetic  
24 state and IPA. This finding aligns with emerging evidence that DNA methylation modulates

25 alternative polyadenylation via CTCF-mediated chromatin looping. Thus, IPAseek provides a  
26 platform to characterize IPA across physiological systems and disease contexts using widely  
27 available bulk RNA-seq data. These IPA events can be further integrated with other regulatory  
28 datasets to elucidate their interplay and functional significance.

## 29 Introduction

30 Cleavage and polyadenylation are a central step in mRNA biogenesis that determines transcript  
31 3' ends and contributes to gene regulatory complexity. Extensive work over the past two decades  
32 has established alternative cleavage and polyadenylation (APA) as a widespread mechanism for  
33 generating transcript isoforms with distinct 3' untranslated regions (3' UTRs)  
34 ([Tian et al. 2005](#); [Wang et al. 2008](#); [Proudfoot 2011](#); [Elkon et al. 2013](#)). Variation in 3' UTR  
35 length can influence post-transcriptional regulation by modulating the presence of *cis*-regulatory  
36 elements that affect mRNA stability, localization, and translation ([Mayr and Bartel 2009](#);  
37 [Berkovits and Mayr 2015](#); [Mayr 2019](#)). Consistent with these regulatory roles, APA has been  
38 implicated in diverse biological processes, including immune activation and tumorigenesis, while  
39 genetic variation near polyadenylation signals (PAS) has been linked to disease-associated traits  
40 ([Sandberg et al. 2008](#); [Neve et al. 2017](#); [Gabel et al. 2024](#)).

41  
42 Intronic polyadenylation (IPA) represents a distinct class of APA in which RNA isoform terminates  
43 within intronic regions, generating truncated mRNA isoforms with altered coding region ([Tian et](#)  
44 [al. 2007](#); [Singh et al. 2018](#)). IPA events can arise through different splicing configurations,  
45 including composite terminal exons or skipped terminal exons, and have been shown to expand  
46 transcript diversity. A well-studied example occurs at the immunoglobulin heavy chain locus,  
47 where developmentally regulated IPA site usage in B cells controls the production of membrane-  
48 bound versus secreted IgM ([Early et al. 1980](#); [Tian et al. 2007](#)). More broadly, IPA has been

49 associated with cell differentiation, immune function, stress responses, and cancer, underscoring  
50 its potential impact on gene regulation and cellular identity ([Ni and Kuperwasser 2016](#); [Singh et  
51 al. 2018](#); [Zhao et al. 2021](#); [Sun et al. 2024](#)). Despite these observations, the factors that shape  
52 cell-type specific IPA site usage remain largely under explored.

53  
54 Although numerous computational approaches have been developed to study APA using RNA  
55 sequencing data, most were designed to detect 3' cleavage events within annotated 3' UTRs and  
56 are not optimized for intronic regions ([Xia et al. 2014](#); [Gruber et al. 2018](#); [Ha et al. 2018](#)).  
57 Detecting IPA poses additional challenges due to lower transcript abundance, complex splicing  
58 patterns, and confounding signals from intron retention. While experimental methods such as 3'-  
59 end sequencing provide high-resolution maps of polyadenylated RNA cleavage sites, their  
60 scalability limits their use in large datasets ([Derti et al. 2012](#); [Lianoglou et al. 2013](#)). Recent bulk  
61 RNA-seq based tools like IPAFinder and InPACT, targeting IPA have begun to address these  
62 limitations but remain sensitive to coverage biases or sequence-based assumptions, highlighting  
63 the need for robust methods specifically tailored to genome-wide analysis of IPA dynamics  
64 ([Scotto-Lavino et al. 2006](#); [Liu et al. 2024](#)).

65  
66 Hematopoietic differentiation provides a well-defined biological framework in which to examine  
67 dynamic RNA processing events across closely related cell states. Epigenetic mechanisms,  
68 including DNA methylation, play essential roles in regulating gene expression during  
69 hematopoiesis and have been implicated in modulating polyadenylation site usage in selective  
70 contexts ([Wood et al. 2008](#); [Cowley et al. 2012](#); [Nanavaty et al. 2020](#)). However, the extent to  
71 which DNA methylation contributes to cell-type specific regulation of IPA remains unclear.

72

73 In this study, we investigate the dynamics of IPA across hematopoietic lineages using genome-  
74 wide transcriptomic and epigenomic data. We develop and apply a computational framework  
75 called IPAseek for systematic identification and quantification of IPA events and use this approach  
76 to characterize lineage-associated patterns of IPA site usage. We further examine the relationship  
77 between DNA methylation and IPA to explore potential epigenetic contributions to cell-type  
78 specific IPA regulation.

## 79 Results

### 80 IPAseek - A Method for IPA Site Detection

81 Changepoints are positions within a data sequence where there is a shift in statistical properties,  
82 such as mean or variance. In the context of bulk RNA-seq data, a used IPA site manifests as a  
83 transition in transcript coverage: there is coverage upstream of the site and a drop or absence of  
84 coverage downstream, creating a clear changepoint signature (composite IPA, **Fig. 1A**). This  
85 pattern can become more complex when splicing events are coupled with 3'-end formation,  
86 resulting in multiple changepoints (e.g., skipped IPA, **Fig. 1D**), and bulk RNA-seq coverage biases  
87 further add to the challenges of changepoint detection ([Killick et al. 2012](#); [Haynes et al. 2017](#)).

88

89 To enable *de novo* detection of IPA sites from individual bulk RNA-seq profiles, we developed a  
90 computationally efficient methodology capable of identifying multiple changepoints within the vast  
91 search space of intronic regions of expressed genes. Our approach leverages the Pruned Exact  
92 Linear Time (PELT) algorithm, a dynamic programming method with linear computational  
93 complexity that applies a penalty for each additional changepoint to balance sensitivity and  
94 specificity ([Killick et al. 2012](#)). Because coverage biases and inherent fluctuations in bulk RNA-  
95 seq data can lead to overfitting (too many changepoints) or underfitting (too few changepoints)  
96 when using a fixed penalty, we integrated PELT with the Changepoints Over a Range of Penalties

97 (CROPS) method, which systematically explores optimal segmentations across a spectrum of  
98 penalty values ([Haynes et al. 2017](#)).

99

100 Our algorithm proceeds as follows: i) Bulk RNA-seq reads are aligned to the genome to generate  
101 coverage profiles; ii) introns of protein-coding genes that exhibit retention are identified and  
102 filtered out to avoid spurious signals; iii) introns of protein-coding genes with contiguous coverage  
103 over at least 100 consecutive bases are selected, as very short covered regions within introns  
104 might not be functionally relevant; iv) PELT is applied to these coverage profiles to detect  
105 changepoints; and v) differential expression upstream and downstream of each detected  
106 changepoint is assessed to confirm a shift in the coverage profile (see Methods and  
107 **Supplementary Fig. S1**). By running changepoint detection over a range of penalties, the  
108 algorithm determines the optimal number of changepoints and evaluates the risk of overfitting at  
109 lower penalties and underfitting at higher penalties (**Fig. 1B-C, 1E-F**). The resulting changepoints  
110 are further evaluated by checking for junction-spanning reads upstream and downstream to  
111 distinguish IPA sites from splicing events.

112

113 A key feature of IPAseek is its use of per-base resolution coverage data from uniquely mapping  
114 reads across introns and flanking coding regions, providing a high-resolution foundation for  
115 changepoint detection. The integration of the PELT algorithm with CROPS, combined with filtering  
116 for intron retention and splicing, enables efficient and robust identification of both single and  
117 multiple changepoints in bulk RNA-seq datasets, thereby enhancing the accuracy and reliability  
118 of IPA site detection (**Fig. 1**).

119 Enhanced Detection of IPA Sites by IPAseek Using Bulk RNA Expression Coverage

120 To evaluate IPAseek performance in *de novo* IPA detection, we applied it to 18 test bulk RNA-  
121 seq samples that also had paired 3'-seq data (see **Supplementary Table S1**) ([Lee et al. 2018](#);

122 [Samur et al. 2018](#)). We first assessed the quality of these bulk RNA-seq samples by comparing  
123 reads aligned to intronic versus exonic regions (**Supplementary Fig. S2A**), and categorized them  
124 as high, moderate, or low-quality based on the relative coverage in these regions.

125  
126 To define ground truth for IPA site validation, we used IPA sites detected with  $\geq 5$  TPM (tags per  
127 million) expression in the 3'-seq atlas as described by Singh et al. ([Lee et al. 2018](#); [Singh et al.](#)  
128 [2018](#)). Currently, 3' end sequencing techniques represent the state-of-the-art methodology for  
129 the precise identification and quantification of transcript 3' ends. Thus, after removal of known  
130 artifacts as described by Singh et al., it was used as a reference standard to evaluate the  
131 performance of IPAseek. Because achieving single-nucleotide precision from bulk RNA-seq is  
132 challenging due to read length and coverage resolution ([Shenker et al. 2015](#); [Arefeen et al.](#)  
133 [2018](#)), IPA sites detected within 350 nt of a ground-truth site were considered true positives (TP),  
134 ground-truth sites missed in RNA-seq were considered false negatives (FN), and IPA sites  
135 detected only in RNA-seq or expressed at  $\leq 5$  TPM in the 3'-seq atlas were categorized as false  
136 positives (FP), hereafter referred to as uncharacterized IPA sites.

137  
138 Using this framework, IPAseek identified IPA events that were not detected by IPAfinder or  
139 InPACT but were supported by ground truth. These included a composite IPA site in *ING5* (106  
140 bp upstream of the 3'-seq sites; **Fig. 2A**) and a skipped IPA site in *EXOC4* (185 bp upstream;  
141 **Fig. 2B**). We validated the 3'-end of the *EXOC4* IPA isoform using 3'-RACE (**Fig. 2C**), confirming  
142 that IPAseek detected a cleavage site 184 bp upstream of the site mapped by 3'-RACE.

143  
144 To further evaluate the biological relevance of IPAseek-detected events, we analyzed PAS  
145 enrichment within  $\pm 200$  nt of IPA sites and observed enrichment for canonical AAUAAA/AUUAAA  
146 motifs across all test samples (**Supplementary Fig. S2B**), consistent with known PAS driven

147 cleavage ([Tian et al. 2005](#)). Approximately 37% of IPaseek-detected IPA isoforms retained  $\leq 25\%$   
148 of the coding sequence (**Supplementary Fig. S2C**), aligning with previous reports that many IPA  
149 isoforms truncate early in the transcription unit ([Singh et al. 2018](#)). Together, these results  
150 demonstrate that IPaseek can detect biologically relevant IPA events.

#### 151 Enhanced IPA Detection with IPaseek

152 We next evaluated IPaseek's performance more systematically by comparing it with two  
153 published IPA detection methods, InPACT and IPAfinder, using precision, recall, and F1-score  
154 metrics (**Fig. 3A**, **Supplementary Fig. S3A**). Across all 18 test samples, IPaseek achieved a  
155 mean precision of 0.16, mean recall of 0.08, and mean F1-score of 0.10, compared with 0.12,  
156 0.06, and 0.08, respectively, for IPAfinder and 0.05, 0.03, and 0.03 for InPACT (see  
157 **Supplementary Table S2**). When stratified by sample quality, IPaseek demonstrated  
158 improvement in metrics within each group; in high-quality samples, precision, recall, and F1-score  
159 reached 0.28, 0.14, and 0.19, respectively, outperforming IPAfinder and InPACT.

160

161 As demonstrated above, under the default settings, IPaseek performs best on high-quality RNA,  
162 with more modest gains over IPAfinder in moderate and low-quality samples. To improve  
163 sensitivity, we relaxed the intron retention filter, defined as the ratio of median intronic to upstream  
164 exon coverage, which increased detectable IPA events but also raised false positives, in a  
165 sample-quality dependent manner. These user-defined thresholds allow IPaseek to be tuned to  
166 different data qualities and tolerances for false positives.

167

168 Even though IPaseek outperforms the other methods, the mean recall remains low, indicating  
169 that a substantial number of IPAs are hard to detect from bulk RNA-seq alone. To examine factors  
170 influencing detection, we compared terminal IPA exon expression between true positives and  
171 false negatives and found that true positives exhibit significantly higher expression than false

172 negatives across multiple test samples (**Fig. 3B, Supplementary Fig. S2D**). This supports IPA  
173 expression level to be a key determinant of detectability in bulk RNA-seq based IPA calling.

174

175 To investigate the validity of false positives, we assessed whether these sites represent authentic  
176 IPA events that are unannotated or missed by our 3'-seq cutoff or instead reflect artifacts. We  
177 examined additional evidence for 3'-end formation by querying PolyASite, RefSeq, Ensembl, and  
178 lowly expressed 3'-seq events ( $\leq 5$  TPM), and by assessing the presence of essential PAS motifs  
179 near the determined cleavage sites ([O'Leary et al. 2016](#); [Herrmann et al. 2020](#); [Martin et al.](#)  
180 [2023](#)). IPAseek detected fewer false positives overall, especially in high-quality samples, than  
181 IPAfinder and InPACT (**Fig. 3C, Supplementary Fig. S3B**). For example, in Test Samples 1, 2,  
182 3, 4, 8, and 11, IPAseek identified 466, 474, 547, 473, 459, and 377 false positives, respectively,  
183 compared with 727, 460, 731, 594, 508, and 547 for IPAfinder and 315, 160, 1874, 1469, 1275,  
184 and 1742 for InPACT.

185

186 Most IPAseek false positives could be cross-referenced to existing databases or supported by  
187 lowly expressed isoforms in the 3'-seq atlas, indicating that many of these sites likely correspond  
188 to bonafide IPA events not captured by the stringent ground-truth definition. False positives from  
189 InPACT and IPAfinder also contained previously unannotated sites, underscoring the potential  
190 of all three methods to expand the landscape of known IPA events. To further evaluate the  
191 biological relevance of unannotated IPAseek false positives, we examined PAS motifs within  $\pm 150$   
192 nt of the potential cleavage sites and observed a strong enrichment for canonical and non-  
193 canonical PAS sequences, with most motifs located close to the IPA sites (**Fig. 3D-E,**  
194 **Supplementary Fig. S3C-D**).

195

196 We conducted further characterization of selected false positive sites and identified a distinct  
197 composite IPA event within the *PARVB* gene (Chr22:44024452-44024923 (+)) and a skipped IPA  
198 event in the *CEACAM8* gene (Chr19:42587622-42588409 (-)) (**Fig. 3F-G**). These IPA events  
199 were not detected by the two other methods, InPACT and IPAfinder. Collectively, these results  
200 establish IPAseek as an effective tool for detecting diverse IPA events across complex  
201 transcriptomic landscapes. Overall, IPAseek represents an advancement in transcriptomics  
202 research with implications for understanding gene expression regulation in both physiological and  
203 disease contexts.

#### 204 IPAseek Reveals Stage-Specific IPA Site Usage During Megakaryocyte Differentiation

205 In this section, we used IPAseek to systematically analyze changes in IPA during the  
206 differentiation of myeloid progenitor K562 cells into megakaryocytes. Bulk RNA-seq data from  
207 eight sequential time points spanning 0-4320 min, with two biological replicates per time point,  
208 were obtained from the GSE213909 dataset ([Bond et al. 2023](#)). IPAseek identified 2,342  
209 confident IPA events across replicates (present in both replicates at any time point), with 225  
210 highly used events ( $\geq 25\%$  usage in  $\geq 2$  samples) selected for detailed analysis.

211  
212 Temporal analysis contrasting early (0-360 min) and late (1440-4320 min) stages revealed distinct  
213 clusters of IPA events with differential usage (**Fig. 4A**). A heatmap visualizing dynamic IPA site  
214 usage showed clear separation of early and late time points, demonstrating progressive rewiring  
215 of the IPA landscape over 72 hours (**Fig. 4A**).

216  
217 Our analysis highlighted two genes, *ACOT8* and *PEX13*, exhibiting contrasting patterns of IPA  
218 site usage linked to differentiation stage (**Fig. 4A**). *ACOT8* showed increased IPA site recognition  
219 at later time points (**Fig. 4B, 4F, Supplementary Fig. S4B**), while *PEX13* displayed reduction in  
220 IPA site usage as differentiation proceeded (**Fig. 4C, 4G, Supplementary Fig. S4C**). Protein-

221 protein interaction network analysis confirmed a direct functional relationship  
222 between *ACOT8* and *PEX13* (**Supplementary Fig. S4A**), supporting coordinated regulation of  
223 peroxisomal function and lipid metabolism during megakaryopoiesis ([Hunt et al. 2012](#); [Lee et al.](#)  
224 [2018](#); [Plessner et al. 2024](#)).

225  
226 To validate these IPA events, we performed 3'-RACE for *ACOT8* and *PEX13*, confirming the IPA  
227 isoforms detected by IPAseek (**Fig. 4D**). 3'-RACE revealed an additional intronic cleavage site  
228 in *PEX13* within the intronic TE, suggesting multiple IPA cleavage sites around this locus.  
229 However, the complexity of the locus and bulk RNA-seq resolution constrained IPAseek's ability  
230 to reliably detect the additional site (detected at multiple time-points but not consistently across  
231 replicates).

232  
233 Megakaryocytic differentiation was induced in K562 cells using 25 nM PMA and confirmed by  
234 phase-contrast microscopy showing hallmark morphological changes (**Fig. 4E**) and upregulation  
235 of *ITGB3* (**Supplementary Fig. S4D**). In-house time-course validation showed *ACOT8* IPA site  
236 usage increasing and *PEX13* IPA site usage decreasing over time (**Fig. 4F, G**), mirroring bulk  
237 RNA-seq patterns determined from the datasets obtained from the public domain  
238 (**Supplementary Fig. S4E-F**).

239  
240 Together, these results demonstrate that IPAseek captures stage-specific IPA events during  
241 megakaryocyte differentiation.

242 Comprehensive Atlas of IPA events Reveals Lineage and Stage Specific Regulation in  
243 Hematopoietic Lineage Cells

244 To systematically characterize IPA events across hematopoietic cell-types, we constructed an  
245 atlas using bulk RNA-seq data from 78 samples representing 10 hematopoietic populations,

246 including differentiated cells (myeloid: granulocytes, monocytes; lymphoid: naïve B cells,  
247 CD4<sup>+</sup>/CD8<sup>+</sup> T cells, NK cells) and undifferentiated precursors (granulocytic, monocytic, erythroid  
248 precursors, and CD34<sup>+</sup> HSPCs).

249

250 To assess whether IPA site usage distinguishes cell lineages and differentiation states, PCA  
251 followed by *k*-means clustering (*k* = 3) of IPA site usage in differentiated cells (GSE184264)  
252 revealed three distinct clusters: granulocytes (myeloid; cluster 1, *n* = 7, silhouette width = 0.36),  
253 lymphoid/monocytes (clusters 2/3, *n* = 13/19, silhouette widths = 0.52/0.56), with overall average  
254 silhouette width of 0.511 (**Fig. 5A, Supplementary Fig. S5C**). To further enhance the separation,  
255 we repeated the PCA and clustering excluding granulocytes (**Supplementary Fig. S5A**). In this  
256 subset, *k*-means (*k*=5) identified five clusters, with most monocyte samples (5 out of 7) forming  
257 Cluster 1, and all lymphoid samples distributed among the remaining clusters. For undifferentiated  
258 precursors (GSE114922), *k*-means clustering (*k* = 4) identified four groups with HSPCs  
259 dominating clusters 1/2 (*n* = 65/17, silhouette widths = 0.49/0.71), overall average silhouette  
260 width 0.52 (**Fig. 5B, Supplementary Fig. S5D**). Given the higher number of HSPC samples, we  
261 performed PCA and clustering after excluding HSPCs (**Supplementary Fig. S5B**). This analysis  
262 identified three groups, with granulocytic precursors and monocytic precursors primarily clustering  
263 together (Cluster 1), demonstrating their lineage similarity. We did this independently for the two  
264 datasets to avoid factors like differences in sample preparation, sequencing techniques etc. that  
265 could convolute the interpretation of results.

266

267 To identify lineage-enriched IPA events, differential analysis between myeloid and lymphoid  
268 lineages identified 981 lineage-enriched IPA events (Wilcoxon rank-sum test, *P* – adj < 0.05),  
269 with hierarchical clustering revealing clear separation (**Fig. 5C**). Similarly, analysis of 36

270 undifferentiated precursor samples identified 372 cell-type specific IPA events (Kruskal-Wallis  
271 test,  $P - \text{adj} < 0.05$ ), showing distinct patterns among granulocytic precursors, monocytic  
272 precursors, and CD34<sup>+</sup> HSPCs (**Fig. 5D**).

273 *RPN1* exemplified lymphoid-enriched IPA site usage (**Fig. 5C, E**), with bulk RNA-seq coverage  
274 confirming ~4.5-fold higher IPA site usage in lymphoid versus myeloid cells (Wilcoxon rank-sum  
275 test,  $P - \text{adj} < 0.05$ ). Together, these validated clusters and differential analyses  
276 demonstrate lineage and stage-specific IPA regulation during hematopoiesis.

277

278 Distinct DNA Methylation Landscapes Surround IPA Sites in Hematopoietic Lineage Cells

279 Although IPA is widespread, gene and cell-type specific, and responsive to stimuli, we know little  
280 about factors enhancing or preventing its usage ([Lee et al. 2018](#); [Singh et al. 2018](#)). Epigenetic  
281 modifications like DNA methylation modulate gene expression without altering DNA sequence,  
282 which is analogous to regulated IPA site recognition, though IPA differs from 3' UTR APA due to  
283 concurrent splicing.

284

285 Prior studies showed DNA methylation ablation increases proximal PAS usage by disrupting  
286 CTCF/cohesin binding ([Nanavaty et al. 2020](#); [Fink et al. 2025](#)), and influences allele-specific IPA  
287 of imprinted genes H13/Nap115 ([Wood et al. 2008](#); [Cowley et al. 2012](#)). No studies have  
288 investigated DNA modifications regulating cell-type specific IPA site usage in its endogenous  
289 state. We therefore examined how endogenous DNA methylation state impacts IPA site selection  
290 across immune cells and cancers.

291

292 To test this systematically, we integrated 198 bulk RNA-seq and 189 BS-seq samples  
293 (RRBS/WGBS) across 12 immune cell-types (**Supplementary Fig. S6A-B**). IPAseek identified

294 IPA sites from RNA-seq data. We compiled IPA signals (PAS) and categorized them as "used"  
295 (within 200 bp of high-confidence IPA: RPKM  $\geq 0.5$ , usage  $\geq 10\%$  in  $\geq 10\%$  samples) or "unused"  
296 to control for sequence biases (**Supplementary Fig. S6C**).

297  
298 Among 910 used and 1,820 unused PAS ( $>5.7\text{M}$  methylation sites), used PAS consistently  
299 exhibited higher surrounding DNA methylation (**Fig. 6A**, KS-test,  $P\text{-value} \leq 9.88 \times 10^{-270}$ ), indicating  
300 methylation influences PAS recognition. This pattern persisted across cell-types: B cells (679 vs.  
301 1358,  $P\text{-value} \leq 5.14 \times 10^{-177}$ ; **Fig. 6B**), CD14<sup>+</sup> monocytes (325 vs. 650,  $P\text{-value} \leq 3.94 \times 10^{-244}$ ); **Fig. 6C**), AML cells (552 vs. 1104,  $P\text{-value} < 10^{-244}$ ; **Fig. 6D**), each analyzing 2.95-3.4M sites.

303  
304 To test if methylation differences extend intron-wide, we performed rolling window (100 bp  
305 windows, 50 bp steps) and equal-tile (500 bins) analyses comparing IPA introns ( $n = 696$ ) to non-  
306 IPA introns from IPA genes ( $n = 1,758$ ) and non-IPA genes ( $n = 1,303$ ) (**Fig. 6E**, **Supplementary**  
307 **Fig. S6D**; KS-test  $P\text{-value} \leq 2 \times 10^{-16}$ ). Higher methylation across IPA introns suggests broader  
308 chromatin context influences cleavage.

309  
310 To determine if methylation changes drive IPA site usage changes, we analyzed GSE184314  
311 (CD4<sup>+</sup> cells), where hypomethylated introns ( $n = 80$ ) were enriched for increased IPA site usage  
312 (Fisher's exact,  $P\text{-value} \leq 2.78 \times 10^{-3}$ ; **Fig. 6G**) while hypermethylated introns ( $n = 135$ ) showed no  
313 association (**Fig. 6F-G**). Reciprocally, low-usage IPA introns ( $n = 19$ ) were enriched for increased  
314 methylation ( $P\text{-value} \leq 5.24 \times 10^{-3}$ ; **Fig. 6H**), but high-usage introns ( $n = 34$ ) showed no association  
315 (**Fig. 6H-I**). Patterns were partially replicated in GSE66117 (**Supplementary Fig. S7A**) but  
316 inconsistent elsewhere (**Supplementary Fig. S7B-F**).

317

318 Collectively, increased intronic DNA methylation correlates with reduced IPA site usage across  
319 contexts, though causality remains untested. Unlike distal CpG island effects on APA, we found  
320 robust proximal ( $\leq 2,500$  bp) methylation-IPA associations (**Fig. 6A-E**).

## 321 Discussion

322 In this study, we present IPaseek, a computational framework that leverages dynamic  
323 programming and changepoint detection algorithms to accurately identify intronic polyadenylation  
324 (IPA) events from bulk RNA-seq data. Unlike previous methods primarily focused on 3' UTR  
325 alternative polyadenylation (APA), IPaseek robustly detects both composite and skipped IPA  
326 isoforms, addressing challenges posed by complex splicing patterns and complexity in intronic  
327 regions. Our benchmarking against paired 3'-end sequencing data and comparison with  
328 established tools such as IPAFinder and InPACT demonstrate that IPaseek achieves superior  
329 sensitivity and precision, enabling the discovery of IPA events which were also validated by  
330 orthogonal techniques like 3'-RACE. This advancement fills a critical in transcriptome analysis,  
331 as IPA has been historically under-characterized despite its emerging importance in gene  
332 regulation ([Gruber et al. 2018](#); [Singh et al. 2018](#)).

333

334 Applying IPaseek to a comprehensive dataset of immune cell-types and hematopoietic  
335 differentiation time courses, we reveal that IPA site usage is highly dynamic and exhibits lineage  
336 and stage-specific patterns. Consistent with previous reports indicating that IPA contributes to  
337 transcriptome diversification in hematopoietic cells, we observed that lymphoid cells exhibit higher  
338 IPA site usage than myeloid cells ([Singh et al. 2018](#)). This finding aligns with the notion that IPA  
339 modulates immune cell identity and function, as exemplified by the well-characterized IPA of the  
340 immunoglobulin M heavy chain (*IGHM*) locus, which controls the production of secreted versus  
341 membrane-bound IgM isoforms in B cells ([Early et al. 1980](#); [Takagaki and Manley 1998](#)). Our

342 temporal profiling during megakaryocyte differentiation further highlights gene-specific IPA  
343 regulation linked to key biological process like autophagy, extending prior observations that IPA  
344 influences gene expression programs during differentiation and stress responses ([Thomas et al.](#)  
345 [2012](#); [Cheng et al. 2020](#)).

346  
347 A major contribution of our study is the integrative analysis of endogenous DNA methylation state  
348 and IPA site usage, which uncovers a previously underappreciated epigenetic layer modulating  
349 IPA. We demonstrate that introns containing IPA sites exhibit elevated DNA methylation  
350 compared to non-IPA introns, a pattern that holds true across multiple cell-types. This observation  
351 is in line with recent findings that DNA methylation can influence APA by modulating chromatin  
352 architecture and the recruitment of RNA processing factors ([Cowley et al. 2012](#); [Smith 2019](#);  
353 [Nanavaty et al. 2020](#); [Jia et al. 2024](#); [Fink et al. 2025](#)). For instance, DNA methylation-  
354 dependent CTCF binding has been shown to regulate APA site choice by altering chromatin  
355 looping and polymerase pausing, thereby affecting transcript isoform diversity. Our results extend  
356 these insights by demonstrating a genome-wide association between methylation and IPA site  
357 usage, supported by enrichment in hypermethylated introns exhibiting reduced IPA site usage.

358  
359 The complex relationship between methylation and IPA site usage we observed suggests that  
360 DNA methylation acts as a context-dependent modulator rather than a simple on/off switch. We  
361 observed that hypermethylation is associated with suppression of IPA. This finding underscores  
362 the need for further mechanistic studies to dissect how epigenetic states integrate with RNA  
363 processing machinery to fine-tune transcript isoform expression in a cell-type and condition-  
364 specific manner.

365

366 The biological implications of our findings are profound. IPA-mediated transcript truncation can  
367 generate non-coding RNAs or protein isoforms with altered functional domains, impacting  
368 processes such as immune signaling, cell fate determination, and disease progression.  
369 Dysregulated IPA has been implicated in cancer, where IPA events can produce truncated tumor  
370 suppressors or oncogenic isoforms, contributing to tumorigenesis and therapy resistance ([Mayr](#)  
371 [and Bartel 2009](#); [Ni and Kuperwasser 2016](#); [Lee et al. 2018](#); [Li et al. 2020](#); [Zhao et al. 2021](#);  
372 [Cheng et al. 2024](#)). Our demonstration that DNA methylation is associated with IPA site selection  
373 suggests an epigenetic mechanism by which cancer cells and immune cells may modulate  
374 transcriptome complexity to adapt to environmental cues or evade immune surveillance.

375

376 IPaseek's scalability and accuracy make it a valuable tool for future investigations into the  
377 regulatory landscape of IPA across diverse biological contexts. Integrating IPA detection with  
378 epigenomic profiling will be important to elucidate how IPA contributes to development and  
379 disease. Future studies leveraging single-cell and long-read sequencing technologies may further  
380 resolve IPA isoform heterogeneity and its functional consequences at higher resolution.  
381 Additionally, experimental perturbation of DNA methylation and chromatin modifiers will be  
382 essential to establish causal links between epigenetic states and IPA regulation.

383

384 In conclusion, our work establishes IPaseek as a platform for dissecting the complexity of IPA  
385 and its epigenetic regulation. By revealing the dynamic interplay between DNA methylation and  
386 IPA site usage, we provide insights into the multi-layered control of gene expression.

387 Methods

388 IPaseek algorithm

389 **Intron Preprocessing**

390 We obtained the human genome annotation (version hg38) from UCSC RefSeq (table name:  
 391 refGene) and flattened the entire genome, annotating each position with a genomic region: intron,  
 392 CDS, 3' UTR, 5' UTR, ncRNA, or intergenic, according to the RefSeq annotation.  
 393 Upstream/downstream regulatory regions were defined as 5 kb extensions from 5' UTR (5' UTR\*)  
 394 and 3' UTR (3' UTR\*) to capture potential regulatory elements influencing the gene expression.  
 395 When annotating overlapping regions, we applied the following priority order: CDS > 5' UTR >  
 396 intron > 3' UTR and CDS+3' UTR. Intronic regions were filtered using:

$$397 \quad L_{intron} \in [500] [150000] \text{ bp}$$

398 where  $L_{intron}$  represents intron length. The lower bound (500 bp) ensures sufficient resolution for  
 399 bulk RNA-seq coverage analysis, while the upper limit (150 kbp) excludes ultra-long introns prone  
 400 to alignment artifacts. Genes lacking coding regions ( $CDS=\emptyset$ ) and introns overlapping  
 401 snoRNA/miRNA loci (UCSC Table Browser, assembly hg38), blacklisted regions (ENCODE), or  
 402 retrotransposons (RepeatMasker) were excluded to focus on protein-coding transcripts and  
 403 minimize confounding signals (Griffiths-Jones 2004; Weber 2005; Griffiths-Jones et al. 2006;  
 404 Griffiths-Jones et al. 2008).

405

## 406 **Sample Preprocessing**

407 For each sample, reads/FASTQ files were aligned to the reference genome using STAR  
 408 (v2.7.10a) with default parameters (Dobin et al. 2013). To reduce false-positive coverage signals,  
 409 only uniquely mapped reads (MAPQ  $\geq$  255) were retained. Gene expression quantification used  
 410 RPKM normalization:

$$411 \quad RPKM = \frac{\text{No. of reads mapped to the gene} \times 10^9}{\text{Total Library Size} \times \text{Gene Length (bp)}}$$

412 Gene expression levels were quantified using RPKM to account for transcript length and  
 413 sequencing depth biases (Mortazavi et al. 2008). Genes with RPKM > 0.5 were considered

414 expressed, a conservative threshold balancing sensitivity and specificity in bulk RNA-Seq  
 415 datasets for IPA site detection.

### 416 **Retained Intron Identification and Filtering**

417 The IPA detection begins with calculating per-base read coverage from uniquely mapping reads  
 418 over filtered introns derived from expressed genes. This step also includes calculating coverage  
 419 for the coding sequence (CDS) flanking the introns under consideration. Following the coverage  
 420 calculations, we focus on removing retained introns, as these are unlikely to contain premature  
 421 cleavage sites (Monteuuis et al. 2019). We established four specific criteria for classifying an  
 422 intron as retained: (1) a minimum of three reads spanning the intron-exon junction, (2) at least  
 423 85% of the intron covered by reads, (3) median coverage over the flanking exons exceeding 0.5  
 424 RPKM, and (4) a ratio of median coverage over the intron to median coverage over the upstream  
 425 exon of at least 5% (Middleton et al. 2017). Any intron meeting all four conditions was classified  
 426 as a retained intron and excluded from further analysis. For the remaining introns, we ensured  
 427 that they have read coverage to contain a potential TE by checking for a read coverage of at least  
 428 five over a contiguous stretch of 100 bp. Introns failing to meet this criterion were also excluded  
 429 from analysis. The next critical step involves detecting significant change points in read coverage  
 430 within the introns, utilizing the Pruned Exact Linear Time (PELT) algorithm.

431

### 432 **Change point Detection in Introns using PELT**

433 The Pruned Exact Linear Time (PELT) algorithm is an efficient and exact method for detecting  
 434 multiple change points in time series data. It minimizes a penalized cost function that balances  
 435 segmentation fit with model complexity, defined as:

$$436 \quad F(n) = \left[ \sum_{i=1}^{k+1} C(y_{(\tau_{i-1}+1):\tau_i}) + \beta k \right]$$

437 where  $C$  is the cost function for each segment,  $k$  is the number of changepoints, and  $\beta$  is a  
438 penalty term that prevents overfitting. The algorithm employs dynamic programming to recursively  
439 calculate the optimal segmentation up to each time point and integrates a pruning step to discard  
440 suboptimal changepoint candidates. The pruning condition ensures computational efficiency by  
441 eliminating candidates that cannot improve the segmentation, reducing the average  
442 computational complexity to  $O(n)$  under certain conditions.

443

444 For bulk RNA-seq analysis, PELT is particularly advantageous due to its ability to adaptively  
445 detect abrupt changes in coverage patterns, which may indicate biologically meaningful events  
446 such as IPA sites. By adjusting parameters like the penalty value

$$447 \quad \beta = 2p \log n$$

448 where  $p$  is the number of variables and  $n$  is the sample size and minimum segment length,  
449 detection sensitivity can be fine-tuned based on dataset characteristics. In this study, we applied  
450 PELT combined with the Changepoints for a Range Of Penalties (CROPS) approach to identify  
451 significant changes in bulk RNA-seq intron coverage. CROPS evaluates multiple penalty values  
452 within a specified range (100 to 10,000 in our case), enabling adaptive sensitivity in changepoint  
453 detection. The implementation was carried out using the `cpt.mean` function from the "changepoint  
454 (version 2.2.1)" R package, which detects mean shifts in time series data. To prepare the data,  
455 intron coverage was extracted from BAM files and filtered to retain regions between 500 and  
456 150,000 base pairs in length. This filtering step reduced computational complexity while focusing  
457 on biologically relevant regions. For antisense strand introns, coverage vectors were reversed to  
458 ensure consistent orientation during processing. The algorithm starts by defining a cost function  
459 that measures how well a segmentation fits the data while penalizing each changepoint to prevent  
460 overfitting. Using dynamic programming, it recursively calculates the optimal segmentation up to  
461 each change point (minimum segment length = 200bp) and applies a pruning rule that eliminates

462 candidate changepoint locations that cannot be part of the optimal solution. This pruning  
463 mechanism significantly reduces computational costs without compromising accuracy.

464

465 By combining PELT with CROPS, we detected abrupt changes in expression coverage patterns  
466 within introns that may correspond to IPA events. The adaptability of this approach allowed us to  
467 optimize detection sensitivity while maintaining computational efficiency critical for genome-scale  
468 analyses. This method demonstrated robust performance in identifying potential IPA sites while  
469 efficiently handling large datasets.

470

### 471 **Changepoint Selection**

472 Next, we assessed whether the detected changepoints met IPA conditions. This was achieved by  
473 calculating the number of spliced reads within  $\pm 50$  bp of each changepoint and evaluating RPKM  
474 coverage within  $\pm 200$  bp. For cases where a single changepoint was identified within an intron,  
475 we checked for Composite IPA conditions, requiring that the median coverage upstream of the  
476 changepoint exceeded the median coverage downstream, with no splice reads present within  $\pm 50$   
477 bp of the changepoint. When multiple changepoints were detected within an intron, adjacent  
478 changepoints were analyzed iteratively in pairs to distinguish between Composite IPA and  
479 Skipped IPA scenarios. For Composite IPA detection using changepoint pairs, two conditions had  
480 to be satisfied: (1) the median coverage upstream of changepoint 1 must exceed the median  
481 coverage downstream of changepoint 1, and (2) the median coverage upstream of changepoint  
482 2 must exceed the median coverage downstream of changepoint 2. Additionally, no splice reads  
483 could be present within  $\pm 50$  bp of either changepoint. For Skipped IPA scenarios, different criteria  
484 were applied: (1) the median coverage upstream of changepoint 1 must be lower than the median  
485 coverage downstream of changepoint 1, while (2) the median coverage upstream of changepoint  
486 2 must exceed the median coverage downstream of changepoint 2. Furthermore, changepoint 1

487 required more than 5% spliced reads within  $\pm 50$  bp, whereas no splice reads could be present  
488 within  $\pm 50$  bp of changepoint 2. (**Supplementary Fig. S1**).

489

490 Once potential IPA sites are identified, we annotate the new Terminal Exon (TE) associated with  
491 each IPA event. For composite IPA events, the end of the upstream coding sequence (CDS) is  
492 designated as the start site of the TE, while the IPA site serves as the end site. For skipped IPA  
493 events, the start of the new TE is defined by changepoint 1, and its end is marked by changepoint  
494 2. This approach ensures accurate annotation of TEs for both types of IPA scenarios.

495

#### 496 **Additional Filtering for Significant IPA Sites**

497 To enhance significance, we applied few additional filters: first, ensuring that coverage (RPKM)  
498 200 bp upstream of each IPA site was greater than 0.5; second, performing differential expression  
499 analysis between three contiguous 100 bp windows upstream and downstream of each  
500 changepoint using DESeq2 (version 1.46.0) ([Love et al. 2014](#)). Only those IPA sites with  
501  $P_{adj} \leq 0.2$  and  $P - \text{value} \leq 0.1$  were deemed significant for downstream analysis. Additionally,  
502 we filtered out new TE for  $\geq 0.5$  RPKM expression. Sites passing these stringent conditions were  
503 classified as valid IPA sites along with their respective annotated TEs designated as TEs for the  
504 corresponding IPA isoforms.

505

#### 506 **Construction of an IPA Atlas and Quantification of IPA Site Usage Across Samples**

507 IPAseek can be utilized to create an atlas of IPA events across multiple samples. Due to the  
508 resolution limitations of bulk RNA-seq, the same IPA site may be detected a few nucleotides apart  
509 in different samples. To address this redundancy when combining IPA events from multiple  
510 samples, we merge IPA events that share the same classification (e.g., composite or skipped)

511 and have TE ends annotated within 100 bp of each other. For these merged IPA events, the TE  
 512 end is assigned as the median of the TE ends from the individual IPA events being combined.

513  
 514 Once the atlas of IPA events is constructed, we assign a binary confidence value to each event.  
 515 IPA events detected in two or more samples are classified as "confident," while those occurring  
 516 in only a single sample are labeled as "not confident."

517 After generating the atlas, we quantify the usage of each IPA event in individual samples. IPAseek  
 518 calculates the relative usage of an IPA isoform by comparing the expression of its TE to the  
 519 combined expression of the gene's last coding exon and its TE. This metric, referred to as IPA  
 520 site usage, provides a normalized measure of how frequently an IPA isoform is utilized relative to  
 521 other isoforms. The formula for IPA site usage is as follows:

$$522 \quad IPA \text{ usage} = \frac{RPKM_{TE}^j}{RPKM_{TE}^j + RPKM_{lastCDS}^j}$$

523 Here  $RPKM_{TE}^j$  and  $RPKM_{lastCDS}^j$  represent the normalized expression levels of the TE determined  
 524 by IPAseek and the last coding exon of the gene, respectively, in each sample  $j$ . This methodology  
 525 allows for comprehensive quantification of IPA isoform usage compared to the full-length isoform  
 526 of the gene and facilitates direct comparison across multiple samples.

### 527 Benchmarking IPAseek Against Established IPA Detection Methods

528 To evaluate the performance of IPAseek in detecting IPA sites, we conducted a comparative  
 529 analysis using four test samples with matched RNA-seq and 3' Seq data (**details provided in**  
 530 **Supplementary\_Table\_S1**). IPA sites expressed at  $\geq 5$  TPM from the 3' Seq data served as the  
 531 ground truth, while corresponding RNA-seq samples were processed using IPAseek to identify  
 532 IPA sites.

533

534 To benchmark IPaseek against existing IPA detection methods, we analyzed the same test  
535 samples using InPACT and IPAFinder. Due to the resolution constraints of RNA-seq data, IPA  
536 sites detected within 350 bp of the ground truth were classified as true positives (TP). Ground  
537 truth IPA sites that were not detected in RNA-seq were considered false negatives (FN), while  
538 sites detected in RNA-seq but absent in the ground truth were categorized as false positives (FP).  
539 These false positives were further referred to as uncharacterized IPA sites in downstream  
540 analyses. This framework enabled a systematic comparison of IPaseek's accuracy and sensitivity  
541 relative to established methods, while accounting for bulk RNA-seq limitations.  
542 To quantify performance, we calculated Precision, Recall, and F1-score for each method. The  
543 formulas used are as follows:

$$544 \quad \textit{Precision} = \frac{TP + FP}{TP}$$

$$545 \quad \textit{Recall} = \frac{TP + FN}{TP}$$

$$546 \quad \textit{F1} = \frac{2 \times (\textit{Precision} \times \textit{Recall})}{\textit{Precision} + \textit{Recall}}$$

547 This comprehensive evaluation allowed us to assess IPaseek's effectiveness in detecting IPA  
548 sites and compare its performance against InPACT and IPAFinder.

#### 549 Annotation Assessment of Uncharacterized IPA Sites

550 All methods identified uncharacterized IPA sites in the test samples that were either absent or not  
551 expressed in the corresponding 3'-seq datasets. To investigate these uncharacterized sites  
552 further, we searched for potential 3'-end annotations in external databases, including PolyAsite,  
553 RefSeq, and Ensembl (O'Leary et al. 2016; Herrmann et al. 2020; Martin et al. 2023). As a first  
554 step, we examined the respective 3'-seq datasets to identify any potential annotations among the  
555 unexpressed candidates. For uncharacterized sites without annotation in the 3'-seq data, we  
556 searched within  $\pm 350$  bp of the annotated 3' UTR end sites in PolyAsite, RefSeq, and Ensembl,

557 following this order of priority. Uncharacterized sites that lacked annotation in any of these  
558 sources were categorized as unannotated.

#### 559 IPA Detection and Quantification During Myeloid Progenitor to Megakaryocyte Differentiation

560 Bulk RNA-seq data from 16 samples, spanning eight time points (0 min, 30 min, 90 min, 180 min,  
561 360 min, 1440 min, 2880 min, and 4320 min) with two replicates per time point, were obtained  
562 from GEO under accession number GSE213909. These samples capture the differentiation of  
563 K562 cells into megakaryocytes. The raw RNA-seq reads were aligned to the human reference  
564 genome (GRCh38) using STAR, and BAM files were filtered to retain only uniquely mapped reads.  
565 IPAseek was then applied to identify IPA sites, construct an atlas of IPA events, and calculate  
566 IPA site usage for each event across all 16 samples.

#### 567 Experimental validation of IPA candidates

#### 568 **3'-RACE**

569 The 3'-RACE was performed using the SMARTer® RACE 5'/3' Kit (Takara, 634858). Briefly, total  
570 cell RNA was extracted from K562 cells and converted to cDNA using SMARTer scribe reverse  
571 transcriptase. The specific genes were amplified from cDNA by gene specific primer in the exonic  
572 region of IPA along UPM followed by a PCR with nested gene specific primer. The amplified  
573 product was then cloned into pRACE vector using infusion cloning and sent to sequencing to  
574 obtain the 3'-end sequence of the RNA. A List of primer used is in the **Supplementary\_Table\_S3**.  
575 Further we generated a FASTA file containing the primer sequences and aligned them to the  
576 reference genome using Bowtie 2. The resulting BAM file was processed with BEDTools to create  
577 a BED file corresponding to primer positions. These coordinates were used to show the position  
578 of RACE primers on the Gviz tracks and have been provided in **Supplementary\_Table\_S3** along  
579 with the primer sequences.

580

#### 581 **Megakaryocyte differentiation from K562 cells**

582 K562 cells were maintained in Iscove's Modified Dulbecco's Medium (IMDM, Gibco™, 12440053)  
583 supplemented with 10% Fetal Bovine serum (FBS, and Penicillin-Streptomycin (Penicillin-  
584 Streptomycin, Gibco™, 15140122). The differentiation of K562 to megakaryocyte was performed  
585 as per Bond et al ([Bond et al. 2023](#)). Briefly,  $1 \times 10^5$  /ml of K562 cells were seeded in complete  
586 media and treated with 25nM of phorbol 12-myristate 13-acetate (PMA) for six spanning time  
587 points (6, 12, 24, 48, 72 and 96 hrs). Media was changed after every 24 hours with  
588 supplementation of fresh PMA. Cells were harvested at 6,12, 24, 48, 72 and 96 hrs. Cells  
589 supplemented with 0.0025% DMSO for 96 hrs were used as control.

590  
591 The differentiation of K562 to megakaryocyte was confirmed by ITGB3, a marker whose  
592 expression increases in megakaryocyte state. We used another gene *KLF1*, erythroid marker as  
593 a negative control whose expression decreases in megakaryocyte state. The expression of the  
594 genes was confirmed using qRT-PCR at 24, 48 and 72 hrs.

595

### 596 **RNA Isolation and real time PCR**

597 Total RNA was isolated using TRIzol reagent (MRC, TR118) by phenol-chloroform extraction  
598 method. Cell pellets were resuspended in TRIzol, phase-separated with chloroform and the RNA  
599 was precipitated with isopropanol followed by washing with 75% ethanol. The pellet was then air  
600 dried and resuspended in RNase-free water. The RNA was converted to cDNA using qScript™  
601 cDNA SuperMix kit (Quantabio, 95048-100) as per manufacture instruction. PowerUp™ SYBR™  
602 Green Master Mix (Applied Biosystems™, A25778) was used to qRT PCR. 25-50 ng RNA  
603 template was used for qRT PCR. Primers used for qRT PCR are listed in  
604 **Supplementary\_Table\_S3**. The relative expression was calculated using  $2^{-\Delta\Delta Ct}$  method and  
605 IPA site usage was calculated by using  $2^{-\Delta\Delta Ct_{IPA}} / (2^{-\Delta\Delta Ct_{IPA}} + 2^{-\Delta\Delta Ct_{FL}})$  ([Singh et](#)  
606 [al. 2018](#)).

## 607 Construction and Analysis of an Immune Cell IPA Atlas

608 To construct an atlas of IPA events in immune cells, we downloaded bulk RNA-seq data from 42  
609 samples under GEO accession number GSE184264 and 36 samples from GSE114922. These  
610 datasets encompass 10 immune cell-types, including naïve B cells, naïve CD4<sup>+</sup> T cells, naïve  
611 CD8<sup>+</sup> T cells, granulocytes, monocytes, natural killer cells, granulocytic precursor cells, monocytic  
612 precursor cells, erythroid precursor cells, and CD34<sup>+</sup> hematopoietic stem progenitor cells (**see**  
613 **Supplementary\_Table\_S4 for details**). Using IPAsseek, we processed the 42 samples to identify  
614 IPA sites and construct the immune cell IPA atlas. Highly confident IPA events were defined as  
615 those with corresponding TEs expressed at RPKM  $\geq 0.5$  in and IPA site usage  $\geq 10\%$  in at least  
616 10% samples of any sample group. Cell-type specific IPA events were identified by selecting  
617 highly confident IPA events detected in at least two samples associated with a given cell-type.  
618 We performed PCA followed by *k*-means clustering on the samples from GSE184264 and  
619 GSE114922 separately to identify the separation between the sample groups based on the IPA  
620 site usage.

621  
622 For differential analysis of IPA site usage between two groups (e.g., cell-types or lineages), we  
623 performed a Wilcoxon rank-sum test for each IPA event and calculated the overall mean IPA site  
624 usage as well as group-specific means (group1\_mean and group2\_mean). *P* – values were  
625 adjusted using the Benjamini-Hochberg FDR correction. IPA events with *P* – adj < 0.05 were  
626 classified as enriched. If group1\_mean was higher, the event was labeled as "group1 enriched";  
627 otherwise, it was labeled as "group2 enriched." For comparisons involving more than two groups  
628 (e.g., multiple cell-types), we applied the Kruskal-Wallis test for each IPA event, followed by FDR  
629 correction. Events with *P* – adj < 0.05 were identified as enriched IPA events.

## 630 Relation between IPA and DNA methylation in immune cells

### 631 Identification of Matched RNA-seq and BS-seq Samples

632 To explore the relationship between IPA site selection and DNA methylation, we curated 198 bulk  
633 RNA-seq and 189 BS-seq samples (RRBS/WGBS) from ENCODE (primary cells), dbGaP  
634 (phs001027), and GEO (GSE165305, GSE214980, GSE184314, GSE156563, GSE130582,  
635 GSE66117, GSE173790, GSE128269, GSE193201). These represented 12 immune cell-types:  
636 B cells, NK cells, T cells, CD14<sup>+</sup> cells, PTCL, CD4<sup>+</sup> cells, B-ALL, T helper cells, HSPCs, PBMC-  
637 reprogrammed iMSCs, AML, and NBM (**see Supplementary\_Table\_S6**).

638

### 639 **IPA Sites Atlas Construction**

640 IPaseek identified IPA events through intron preprocessing, STAR alignment with RPKM  
641 normalization, PELT changepoint detection, and validation using splice read exclusion and RPKM  
642 thresholds. Events within 100 bp were merged, and confidence was assigned based on detection  
643 frequency.

644

### 645 **Methylation Data Processing Pipeline**

646 BS-seq reads were trimmed (trim-galore v0.6.10), aligned to bisulfite-converted GRCh38  
647 (Bismark v0.24.1), deduplicated, and methylation calls extracted ([Krueger and Andrews 2011](#)).

648

### 649 **Methylation Atlas Generation**

650 Using edgeR v4.4.2 ([Robinson et al. 2010](#); [Chen et al. 2018](#)), sites were filtered ( $\geq 10$  reads,  
651  $\geq 50\%$  methylation in  $\geq 2$  samples) and annotated to GRCh38 gene bodies. Methylation levels were  
652 calculated as proportion of methylated reads. Next we the PAS sites into used and unused and  
653 performed the Integrative Analysis of Methylation Patterns around IPA Sites using Rolling Window  
654 and Equal-Tile Analysis (Refer to Supplementary Methods for details)

655 Data sets

656 A detailed description of datasets used in study is provided in **Supplementary\_Table\_S1, S4,**  
657 **S5&S6**. 3'-seq and bulk RNA-seq datasets for testing and benchmarking IPaseek were  
658 downloaded from GEO accession number GSE111793 & GSE111310 (Test Sample 2-18, Test  
659 Sample 1&2 provided as raw data on IPaseek GitHub repository). Bulk RNA-seq datasets for  
660 studying the temporal dynamics of IPA during megakaryocyte differentiation was downloaded  
661 from GEO accession number GSE213909. Bulk RNA-seq datasets from GSE184264 and  
662 GSE114922 were used for making IPA atlas for immune cells. For understanding the interplay  
663 between IPA site usage and DNA methylation, the RNA-seq and matched BS-seq (RRBS or  
664 WGBS) samples were used from ENCODE, GEO (GSE165305, GSE214980, GSE184314,  
665 GSE156563, GSE130582, GSE66117, GSE173790, GSE128269, GSE193201), dbGaP  
666 (phs001027).

667 Software availability

668 IPaseek is an open-source method available on GitHub repository  
669 (<https://github.com/isinghlab/IPaseek.git>) and as Supplemental Code.

670 Competing interest statement

671 The authors declare that they have no competing interests.

672 Acknowledgements

673 This work was supported by 1R21NS121945, 1R01CA282251, CPRIT RP230204, and Texas  
674 A&M Health Science Center Seedling grant awarded to IS. We thank all Singh Lab members in  
675 the group for their assistance and constructive suggestions. We also gratefully acknowledge the  
676 technical support by the High-Performance Research Computing resources at Texas A&M

677 University. Artificial intelligence-based tools were used exclusively for language editing and text  
678 refinement.

679 Author contributions

## 680 **Authors and Affiliations**

681 Department of Cell Biology & Genetics, Texas A&M University Health Science Center, Bryan,  
682 TX

683 Richa Rashmi, Pranita Borkar, Sumana Mallick, Taylor Hubbs, Ari Aviles, Daniel Chung and  
684 Irtisha Singh

685

686 Department of Biomedical Engineering, Texas A&M University, College Station, TX

687 Abhinaya Muruganandham, Sumana Mallick, Irtisha Singh

688

689 Interdisciplinary Program in Genetics and Genomics, Texas A&M University, College Station,  
690 TX

691 Ari Aviles, Irtisha Singh

692

## 693 **Contributions**

694 IS and RR conceived and designed the study. IS supervised the study and data analysis. RR  
695 implemented the idea, wrote the original code, and analyzed the data. AM helped in data analysis.

696 PB and SM performed the 3'-RACE experiments and performed the K562 differentiation  
697 experiments respectively and analyzed the outcoming data. AA and DC performed the RNA

698 isolation and cDNA conversion for the K562 differentiation study. TH performed the qRT-PCR for  
699 the K562 differentiation study. RR and IS wrote the manuscript with input from all authors. All

700 authors approved the final version submitted.

701

702 **Corresponding author**

703 Correspondence to Irtisha Singh.

704 **References**  
705

706 **Uncategorized References**

707 Arefeen A, Liu J, Xiao X, Jiang T. 2018. TAPAS: tool for alternative polyadenylation site analysis.

708 *Bioinformatics* **34**: 2521–2529. doi:10.1093/bioinformatics/bty110.

709 Berkovits BD, Mayr C. 2015. Alternative 3' UTRs act as scaffolds to regulate membrane protein

710 localization. *Nature* **522**: 363–367. doi:10.1038/nature14321.

711 Bond ML, Davis ES, Quiroga IY, Dey A, Kiran M, Love MI, Won H, Phanstiel DH. 2023. Chromatin

712 loop dynamics during cellular differentiation are associated with changes to both anchor

713 and internal regulatory features. *Genome Res* **33**: 1258–1268.

714 doi:10.1101/gr.277397.122.

715 Chen Y, Pal B, Visvader JE, Smyth GK. 2018. Differential methylation analysis of reduced

716 representation bisulfite sequencing experiments using edgeR. *F1000Research* **6**

717 doi:10.12688/f1000research.13196.2.

718 Cheng LC, Zheng D, Baljinnyam E, Sun F, Ogami K, Yeung PL, Hoque M, Lu C-W, Manley JL,

719 Tian B et al. 2020. Widespread transcript shortening through alternative polyadenylation

720 in secretory cell differentiation. *Nature Communications* **2020** *11:1* **11**

721 doi:10.1038/s41467-020-16959-2.

722 Cheng X, Jiang G, Zhou X, Wang J, Zhao Z, Zhang J, Ni T. 2024. The landscape and clinical

723 relevance of intronic polyadenylation in human cancers. *J Genet Genomics* **51**: 1030–

724 1039. doi:10.1016/j.jgg.2024.04.014.

- 725 Cowley M, Wood AJ, Bohm S, Schulz R, Oakey RJ. 2012. Epigenetic control of alternative mRNA  
726 processing at the imprinted Herc3/Nap115 locus. *Nucleic Acids Res* **40**: 8917–8926.  
727 doi:10.1093/nar/gks654.
- 728 Derti A, Garrett-Engele P, Macisaac KD, Stevens RC, Sriram S, Chen R, Rohl CA, Johnson JM,  
729 Babak T. 2012. A quantitative atlas of polyadenylation in five mammals. *Genome Res* **22**:  
730 1173–1183. doi:10.1101/gr.132563.111.
- 731 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras  
732 TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15–21.  
733 doi:10.1093/bioinformatics/bts635.
- 734 Early P, Rogers J, Davis M, Calame K, Bond M, Wall R, Hood L. 1980. Two mRNAs can be  
735 produced from a single immunoglobulin mu gene by alternative RNA processing  
736 pathways. *Cell* **20**: 313–319. doi:10.1016/0092-8674(80)90617-0.
- 737 Elkon R, Ugalde AP, Agami R. 2013. Alternative cleavage and polyadenylation: extent, regulation  
738 and function. *Nat Rev Genet* **14**: 496–506. doi:10.1038/nrg3482.
- 739 Fink EE, Zhang Y, Santo B, Siddavatam A, Ou R, Nanavaty V, Lee BH, Ting AH. 2025. Heat  
740 shock induces alternative polyadenylation through dynamic DNA methylation and  
741 chromatin looping. *Cell Stress Chaperones* **30**: 100084.  
742 doi:10.1016/j.cstres.2025.100084.
- 743 Gabel AM, Belleville AE, Thomas JD, McKellar SA, Nicholas TR, Banjo T, Crosse EI, Bradley RK.  
744 2024. Multiplexed screening reveals how cancer-specific alternative polyadenylation  
745 shapes tumor growth in vivo. *Nat Commun* **15**: 959. doi:10.1038/s41467-024-44931-x.

- 746 Griffiths-Jones S. 2004. The microRNA Registry. *Nucleic Acids Res* **32**: D109–111.  
747 doi:10.1093/nar/gkh023.
- 748 Griffiths-Jones S, Grocock RJ, van Dongen S, Bateman A, Enright AJ. 2006. miRBase: microRNA  
749 sequences, targets and gene nomenclature. *Nucleic Acids Res* **34**: D140–144.  
750 doi:10.1093/nar/gkj112.
- 751 Griffiths-Jones S, Saini HK, van Dongen S, Enright AJ. 2008. miRBase: tools for microRNA  
752 genomics. *Nucleic Acids Res* **36**: D154–158. doi:10.1093/nar/gkm952.
- 753 Gruber AJ, Schmidt R, Ghosh S, Martin G, Gruber AR, van Nimwegen E, Zavolan M. 2018.  
754 Discovery of physiological and cancer-related regulators of 3' UTR processing with  
755 KAPAC. *Genome Biol* **19**: 44. doi:10.1186/s13059-018-1415-3.
- 756 Ha KCH, Blencowe BJ, Morris Q. 2018. QAPA: a new method for the systematic analysis of  
757 alternative polyadenylation from RNA-seq data. *Genome Biol* **19**: 45. doi:10.1186/s13059-  
758 018-1414-4.
- 759 Haynes K, Eckley IA, Fearnhead P. 2017. Computationally Efficient Change-point Detection for a  
760 Range of Penalties. *Journal of Computational and Graphical Statistics* **26**: 134–143.  
761 doi:10.1080/10618600.2015.1116445.
- 762 Herrmann CJ, Schmidt R, Kanitz A, Artimo P, Gruber AJ, Zavolan M. 2020. PolyASite 2.0: a  
763 consolidated atlas of polyadenylation sites from 3' end sequencing. *Nucleic Acids Res* **48**:  
764 D174–D179. doi:10.1093/nar/gkz918.
- 765 Hunt MC, Siponen MI, Alexson SE. 2012. The emerging role of acyl-CoA thioesterases and  
766 acyltransferases in regulating peroxisomal lipid metabolism. *Biochim Biophys Acta* **1822**:  
767 1397–1410. doi:10.1016/j.bbadis.2012.03.009.

- 768 Jia J, Fan H, Wan X, Fang Y, Li Z, Tang Y, Zhang Y, Huang J, Fang D. 2024. FUS reads histone  
769 H3K36me3 to regulate alternative polyadenylation. *Nucleic Acids Res* **52**: 5549–5571.  
770 doi:10.1093/nar/gkae184.
- 771 Killick R, Fearnhead P, Eckley IA. 2012. Optimal Detection of Changepoints With a Linear  
772 Computational Cost. *Journal of the American Statistical Association* **107**: 1590–1598.  
773 doi:10.1080/01621459.2012.737745.
- 774 Krueger F, Andrews SR. 2011. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq  
775 applications. *Bioinformatics* **27**: 1571–1572. doi:10.1093/bioinformatics/btr167.
- 776 Lee SH, Singh I, Tisdale S, Abdel-Wahab O, Leslie CS, Mayr C. 2018. Widespread intronic  
777 polyadenylation inactivates tumour suppressor genes in leukaemia. *Nature* **561**: 127–131.  
778 doi:10.1038/s41586-018-0465-8.
- 779 Li Y, Schaefer B, Zou X, Zhang M, Heyd F, Sun W, Zhang B, Li G, Liang W, He Y et al. 2020.  
780 Pan-tissue analysis of allelic alternative polyadenylation suggests widespread functional  
781 regulation. *Mol Syst Biol* **16**: e9367. doi:10.15252/msb.20199367.
- 782 Lianoglou S, Garg V, Yang JL, Leslie CS, Mayr C. 2013. Ubiquitously transcribed genes use  
783 alternative polyadenylation to achieve tissue-specific expression. *Genes Dev* **27**: 2380–  
784 2396. doi:10.1101/gad.229328.113.
- 785 Liu X, Chen H, Li Z, Yang X, Jin W, Wang Y, Zheng J, Li L, Xuan C, Yuan J et al. 2024. InPACT:  
786 a computational method for accurate characterization of intronic polyadenylation from  
787 RNA sequencing data. *Nat Commun* **15**: 2583. doi:10.1038/s41467-024-46875-8.
- 788 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-  
789 seq data with DESeq2. *Genome Biol* **15**: 550. doi:10.1186/s13059-014-0550-8.

- 790 Martin FJ, Amode MR, Aneja A, Austine-Orimoloye O, Azov AG, Barnes I, Becker A, Bennett R,  
791 Berry A, Bhai J et al. 2023. Ensembl 2023. *Nucleic Acids Res* **51**: D933–D941.  
792 doi:10.1093/nar/gkac958.
- 793 Mayr C. 2019. What Are 3' UTRs Doing? *Cold Spring Harb Perspect Biol* **11**  
794 doi:10.1101/cshperspect.a034728.
- 795 Mayr C, Bartel DP. 2009. Widespread shortening of 3'UTRs by alternative cleavage and  
796 polyadenylation activates oncogenes in cancer cells. *Cell* **138**: 673–684.  
797 doi:10.1016/j.cell.2009.06.016.
- 798 Middleton R, Gao D, Thomas A, Singh B, Au A, Wong JJ, Bomane A, Cosson B, Eyraas E, Rasko  
799 JE et al. 2017. IRFinder: assessing the impact of intron retention on mammalian gene  
800 expression. *Genome Biol* **18**: 51. doi:10.1186/s13059-017-1184-4.
- 801 Monteuis G, Wong JLL, Bailey CG, Schmitz U, Rasko JEJ. 2019. The changing paradigm of  
802 intron retention: regulation, ramifications and recipes. *Nucleic Acids Res* **47**: 11497–  
803 11513. doi:10.1093/nar/gkz1068.
- 804 Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. 2008. Mapping and quantifying  
805 mammalian transcriptomes by RNA-Seq. *Nat Methods* **5**: 621–628.  
806 doi:10.1038/nmeth.1226.
- 807 Nanavaty V, Abrash EW, Hong C, Park S, Fink EE, Li Z, Sweet TJ, Bhasin JM, Singuri S, Lee BH  
808 et al. 2020. DNA Methylation Regulates Alternative Polyadenylation via CTCF and the  
809 Cohesin Complex. *Mol Cell* **78**: 752–764 e756. doi:10.1016/j.molcel.2020.03.024.

- 810 Neve J, Patel R, Wang Z, Louey A, Furger AM. 2017. Cleavage and polyadenylation: Ending the  
811 message expands gene regulation. *RNA Biol* 14: 865–890.  
812 doi:10.1080/15476286.2017.1306171.
- 813 Ni TK, Kuperwasser C. 2016. Premature polyadenylation of MAGI3 produces a dominantly-acting  
814 oncogene in human breast cancer. *Elife* 5 doi:10.7554/eLife.14730.
- 815 O'Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, Rajput B, Robbertse B,  
816 Smith-White B, Ako-Adjei D et al. 2016. Reference sequence (RefSeq) database at NCBI:  
817 current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**:  
818 D733–745. doi:10.1093/nar/gkv1189.
- 819 Plessner M, Thiele L, Hofhuis J, Thoms S. 2024. Tissue-specific roles of peroxisomes revealed  
820 by expression meta-analysis. *Biol Direct* **19**: 14. doi:10.1186/s13062-024-00458-1.
- 821 Proudfoot NJ. 2011. Ending the message: poly(A) signals then and now. *Genes Dev* 25: 1770–  
822 1782. doi:10.1101/gad.17268411.
- 823 Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential  
824 expression analysis of digital gene expression data. *Bioinformatics* 26: 139–140.  
825 doi:10.1093/bioinformatics/btp616.
- 826 Samur MK, Minvielle S, Gulla A, Fulciniti M, Cleyngen A, Aktas Samur A, Szalat R, Shamma M,  
827 Magrangeas F, Tai YT et al. 2018. Long intergenic non-coding RNAs have an independent  
828 impact on survival in multiple myeloma. *Leukemia* 32: 2626–2635. doi:10.1038/s41375-  
829 018-0116-y.

- 830 Sandberg R, Neilson JR, Sarma A, Sharp PA, Burge CB. 2008. Proliferating cells express mRNAs  
831 with shortened 3' untranslated regions and fewer microRNA target sites. *Science* **320**:  
832 1643–1647. doi:10.1126/science.1155390.
- 833 Scotto-Lavino E, Du G, Frohman MA. 2006. 3' end cDNA amplification using classic RACE. *Nat*  
834 *Protoc* **1**: 2742–2745. doi:10.1038/nprot.2006.481.
- 835 Shenker S, Miura P, Sanfilippo P, Lai EC. 2015. IsoSCM: improved and alternative 3' UTR  
836 annotation using multiple change-point inference. *RNA* **21**: 14–27.  
837 doi:10.1261/rna.046037.114.
- 838 Singh I, Lee SH, Sperling AS, Samur MK, Tai YT, Fulciniti M, Munshi NC, Mayr C, Leslie CS.  
839 2018. Widespread intronic polyadenylation diversifies immune cell transcriptomes. *Nat*  
840 *Commun* **9**: 1716. doi:10.1038/s41467-018-04112-z.
- 841 Smith LM. 2019. Epigenetic Regulation of mRNA Polyadenylation Site Selection. *Plant Physiol*  
842 **180**: 7–9. doi:10.1104/pp.19.00374.
- 843 Sun J, Kim J-Y, Jun S, Park M, de Jong E, Chang J-W, Cheng S, Fan D, Chen Y, Griffin TJ et al.  
844 2024. Dichotomous intronic polyadenylation profiles reveal multifaceted gene functions in  
845 the pan-cancer transcriptome. *Experimental & Molecular Medicine* **2024** *56*:10 **56**  
846 doi:10.1038/s12276-024-01289-w.
- 847 Takagaki Y, Manley JL. 1998. Levels of Polyadenylation Factor CstF-64 Control IgM Heavy Chain  
848 mRNA Accumulation and Other Events Associated with B Cell Differentiation. *Molecular*  
849 *Cell* **2** doi:10.1016/S1097-2765(00)80291-9.

- 850 Thomas PE, Wu X, Liu M, Gaffney B, Ji G, Li QQ, Hunt AG. 2012. Genome-Wide Control of  
851 Polyadenylation Site Choice by CPSF30 in Arabidopsis. *The Plant Cell* **24**  
852 doi:10.1105/tpc.112.096107.
- 853 Tian B, Hu J, Zhang H, Lutz CS. 2005. A large-scale analysis of mRNA polyadenylation of human  
854 and mouse genes. *Nucleic Acids Res* **33**: 201–212. doi:10.1093/nar/gki158.
- 855 Tian B, Pan Z, Lee JY. 2007. Widespread mRNA polyadenylation events in introns indicate  
856 dynamic interplay between polyadenylation and splicing. *Genome Res* **17**: 156–165.  
857 doi:10.1101/gr.5532707.
- 858 Wang ET, Sandberg R, Luo S, Khrebtkova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP,  
859 Burge CB. 2008. Alternative isoform regulation in human tissue transcriptomes. *Nature*  
860 **456**: 470–476. doi:10.1038/nature07509.
- 861 Weber MJ. 2005. New human and mouse microRNA genes found by homology search. *FEBS J*  
862 **272**: 59–73. doi:10.1111/j.1432-1033.2004.04389.x.
- 863 Wood AJ, Schulz R, Woodfine K, Koltowska K, Beechey CV, Peters J, Bourc'his D, Oakey RJ.  
864 2008. Regulation of alternative polyadenylation by genomic imprinting. *Genes Dev* **22**:  
865 1141–1146. doi:10.1101/gad.473408.
- 866 Xia Z, Donehower LA, Cooper TA, Neilson JR, Wheeler DA, Wagner EJ, Li W. 2014. Dynamic  
867 analyses of alternative polyadenylation from RNA-seq reveal a 3'-UTR landscape across  
868 seven tumour types. *Nat Commun* **5**: 5274. doi:10.1038/ncomms6274.
- 869 Zhao Z, Xu Q, Wei R, Wang W, Ding D, Yang Y, Yao J, Zhang L, Hu YQ, Wei G et al. 2021.  
870 Cancer-associated dynamics and potential regulators of intronic polyadenylation revealed

871 by IPAfinder using standard RNA-seq data. *Genome Res* **31**: 2095–2106.  
872 doi:10.1101/gr.271627.120.

873

874 **Fig. 1: Changepoint Detection of IPA Sites by PELT.**

- 875 A. Top - Bulk RNA-seq coverage (raw read count of uniquely mapping reads) over a  
876 genomic locus annotated as intron by RefSeq. The sashimi plot shows the splicing  
877 complexity at the locus. Bottom - Peaks (raw read count of mapped reads at the locus)  
878 detected by 3'-seq, a high throughput 3'-end detection and quantification approach. This  
879 is a representation of a composite IPA event. The gene locus presented here is located  
880 on the sense strand of DNA.
- 881 B. Elbow plot demonstrates the number of changepoints (x-axis) detected by PELT over a  
882 range of penalties (y-axis) over the locus shown in A. This plot shows detection of one  
883 change point generated due to 3'-end formation.
- 884 C. Plot shows the read coverage over the intron in A. It highlights the two segments  
885 detected by PELT with the maximum difference in mean of the segments.
- 886 D. As A but shows the locus of a skipped IPA event. The gene locus presented here is  
887 located on the sense strand of DNA.
- 888 E. As B but shows detection of two change points for the locus shown in D. One change  
889 point explains the splicing and the second change point highlights the 3'-end formation.
- 890 F. As in C but for the locus shown in D. PELT detects three segments with different means  
891 using the read coverage over the intron shown in D.

892 **Fig. 2: Comprehensive analysis of IPA events detected by IPAseek.**

- 893 A. Same as 1A but a composite IPA event only detected by IPAseek. The highlighted  
894 region shows the new terminal exon (TE), where the termination site was detected as a

895 change point by IPaseek. The gene locus presented here is located on the sense strand  
896 of DNA.

897 B. Similar visualization as in A but depicting a skipped IPA event detected only by IPaseek.

898 Top track here indicates the position of the two nested primers used for 3'-RACE  
899 validation. The blue highlighted region shows the TE detected by IPaseek and the pink  
900 highlighted region shows the cleavage site detected by 3'-RACE ( $\pm 20$  nt). The gene  
901 locus presented here is located on the sense strand of DNA.

902 C. Validation of 3'-end detected by IPaseek in B using 3'-RACE. Total RNA extracted from  
903 K562 cells were used to perform 3'-RACE.

904 **Fig. 3: Comparative Analysis of IPaseek Detected IPA Events with Existing Methods.**

905 A. A comprehensive evaluation of the performance metrics for IPaseek, InPACT, and  
906 IPAFinder in detecting IPA events across multiple 6 test samples (ordered by sample  
907 quality). The analysis employs three key metrics: Precision, Recall, and F1-Score. The  
908 IPA 3'-ends detected by 3'-seq (expressed at  $\geq 5$  TPM) were utilized as ground truth.  
909 Precision measures the proportion of correctly identified IPA events among all detected  
910 events, while recall quantifies the fraction of true IPA events successfully identified by  
911 each method. The F1-Score, the harmonic mean of precision and recall, provides a  
912 balanced measure of overall performance.

913 B. Comparison of bulk RNA-seq expression levels of TEs, represented as  $\log_2(\text{RPKM} + 1)$ ,  
914 between IPA events detected by both IPaseek and 3'-seq (True Positives) and those  
915 identified exclusively by 3'-seq (False Negatives). Statistical significance of the  
916 difference in TE expression between True Positives and False Negatives was assessed  
917 using the Kolmogorov-Smirnov test ( $P$  – value  $< 0.001$  of Test sample 1,2,4,8 and  $P$  –  
918 value  $< 0.01$  of Test sample 14).

- 919 C. This figure presents a quantitative comparison of previously uncharacterized IPA events  
 920 (False Positives) identified by IPAseek, InPACT, and IPAFinder across the 6 test  
 921 Samples (ordered by sample quality). Previously uncharacterized IPAs are defined as  
 922 those detected using RNA-seq coverage profiles that either are not present in the 3'-seq  
 923 atlas or have low expression ( $\leq 5$  TPM) in the 3'-seq atlas. The color-coded segments  
 924 within each bar represent the distribution of annotation sources for these  
 925 uncharacterized IPA sites, including the 3'-seq atlas (expressed  $\leq 5$  TPM), PolyASite  
 926 database, RefSeq, and Ensembl annotations, prioritized in that order. The grey region  
 927 indicates IPA events detected by each method that lack any prior annotation.
- 928 D. Proportion of uncharacterized IPA events without annotations detected in 6 Test  
 929 Samples (ordered by sample quality, grey events in C for IPAseek) with a canonical/non-  
 930 canonical PAS within  $\pm 150$ nt of the detected IPA site.
- 931 E. Density plot of the distance of the nearest PAS from the IPA cleavage site detected by  
 932 IPAseek.
- 933 F. Same as 1A. Shows an uncharacterized composite IPA isoform without annotation  
 934 detected by IPAseek but not by InPACT and IPAFinder. The gene locus presented here  
 935 is located on the sense strand of DNA.
- 936 G. Same as 1A. Shows an uncharacterized skipped IPA isoform without annotation  
 937 detected by IPAseek but not by InPACT and IPAFinder. The gene locus presented here  
 938 is located on the antisense strand of DNA.

939 **Fig. 4: Temporal Analysis of IPA During Myeloid Progenitor to Megakaryocyte**

940 **Differentiation**

- 941 A. IPA sites with differential usage ( $\geq 10\%$  usage difference) with respect to full-length  
 942 mRNA in myeloid progenitor cells (K562) to megakaryocytic differentiation ( $n = 166$ ). The  
 943 sites were determined by the calculating the difference in IPA site usage of all the IPA

944 events, using the two early and two end differentiation time points. Each row is a unique  
945 IPA isoform, while columns are samples captured at different time points of the  
946 differentiation process. The color scheme indicates regulation status, with higher usage  
947 shown in orange and lower usage in blue. IPA site usage is quantified on a 0-1 scale  
948 relative to full-length isoform usage, with color intensity reflecting the degree of IPA site  
949 usage.

950 B. RNA-seq read coverage (in TPM) over *ACOT8* locus across different timepoints of  
951 myeloid progenitor (K562) to megakaryocytic differentiation. The read coverage was  
952 determined using uniquely mapping reads. Topmost track indicates the position of the  
953 two nested primers used for 3'-RACE validation. The blue highlighted region shows the  
954 TE determined by IPaseek while pink highlighted region shows the cleavage site  
955 detected by 3'-RACE ( $\pm 20$  nt). The gene locus presented here is located on the  
956 antisense strand of DNA.

957 C. Like B but showing *PEX13* locus. The gene locus presented here is located on the  
958 sense strand of DNA.

959 D. Validation of 3'-end detected by IPaseek in B & C using 3'-RACE. Total cell RNA  
960 extracted from K562 cells were used to perform 3'-RACE.

961 E. Morphological assessment of myeloid progenitor cells. Phase contrast microscopy (20 $\times$   
962 magnification; BioTek Lionheart FX Automated Microscope) captured progressive  
963 morphological changes during PMA-induced megakaryocytic differentiation of myeloid  
964 progenitor cells. Columns represent treatment groups: untreated control (left), 0.0025%  
965 DMSO vehicle (middle), and 25 nM PMA-treated cells (right). Rows correspond to time  
966 points: 24 h (top), 48 h (middle), and 72 h (bottom). PMA-treated cells exhibited hallmark  
967 differentiation features including increased cell size, cytoplasmic expansion, and  
968 enhanced substrate adhesion compared to controls. Scale bars: 10  $\mu$ m.

- 969 F. IPA isoform usage in *ACOT8* during myeloid progenitor differentiation to  
970 megakaryocytes. The left panel shows IPA site usage measured across eight time points  
971 (0 min to 4320 min) following PMA treatment, quantified using IPAseek analysis of bulk  
972 RNA-seq data from GSE213909. The right panel presents IPA isoform usage, identified  
973 by 3'-RACE, measured at six time points (6 h to 96 h) post-PMA treatment or vehicle  
974 control (DMSO 96 h), validated by qRT-PCR. Together, these plots demonstrate an  
975 increase in *ACOT8* IPA isoform usage throughout differentiation.
- 976 G. Same as F but for *PEX13*. It demonstrates a decrease in *PEX13* IPA isoform usage  
977 throughout differentiation.

978 **Fig. 5: Immune cell IPA atlas and functional enrichment analysis.**

- 979 A. Principal Component Analysis (PCA) of differentiated immune cell-types from the GEO  
980 dataset GSE184264 based on IPA site usage. Three clusters identified through *k*-means  
981 clustering are indicated by dotted enclosures. The plot reveals clear segregation  
982 between myeloid (Clusters 1 and 2) and lymphoid (Clusters 2 and 3) lineages, reflecting  
983 distinct lineage-specific IPA patterns. The analysis includes 4,994 IPA events across 39  
984 samples.
- 985 B. PCA of undifferentiated immune cell-types from GEO dataset GSE114922 based on IPA  
986 site usage. Four clusters identified through *k*-means clustering are indicated by dotted  
987 enclosures. The analysis reveals separation between individual cell-types, emphasizing  
988 their unique IPA signatures (No. of IPA = 592, No. of Samples = 108).
- 989 C. Heatmap of IPA site usage across differentiated immune cell-types, showcasing myeloid  
990 and lymphoid specific enrichment patterns. Each row represents a unique IPA isoform (*n*  
991 = 1,863), while columns correspond to individual immune cell samples (*n* = 42). Rows  
992 are color-coded to indicate enrichment status: myeloid enriched IPAs (seafoam green)  
993 and lymphoid enriched IPAs (dark orange). IPA site usage is quantified on a 0-1 scale

- 994 relative to full-length isoform usage, with color intensity reflecting the degree of IPA site  
 995 usage. Statistical significance was determined using the Wilcoxon rank-sum test with  
 996 Benjamini-Hochberg FDR correction ( $P - \text{adj} < 0.05$ ).
- 997 D. Heatmap of IPA site usage across undifferentiated immune cell-types, focusing on the  
 998 cell-types shown in (B). Each row represents a unique IPA isoform ( $n = 431$ ), and  
 999 columns denote individual immune cell samples ( $n = 108$ ). Statistical significance for  
 1000 enrichment in specific cell-types was assessed using the Kruskal-Wallis test with  
 1001 Benjamini-Hochberg FDR correction ( $P - \text{adj} < 0.05$ ). IPA site usage is displayed on a  
 1002 0-1 scale, with color intensity indicating the degree of IPA site usage.
- 1003 E. Differential IPA site usage between myeloid and lymphoid lineages, with *RPN1* shown  
 1004 as an example of a lineage-specific IPA event. Bulk RNA-seq read coverage (in TPM)  
 1005 over the *RPN1* locus is visualized across different immune cell-types. Read coverage  
 1006 was calculated using uniquely mapping reads.

1007 **Fig. 6: Epigenetic Regulation of IPA in Immune Cells.**

- 1008 A. Integrated Methylation Analysis at IPA Sites. Metagene analysis of DNA methylation  
 1009 patterns within  $\pm 2500$  bp of PAS revealed distinct epigenetic landscapes between used  
 1010 (orange) and unused (navy) PAS loci. Used PAS sites exhibited significantly higher  
 1011 methylation levels across all samples compared to unused PAS (Kolmogorov-Smirnov  
 1012 test,  $P - \text{value} \leq 9.88 \times 10^{-270}$ ), with mean methylation differences persisting across the  
 1013 entire genomic window. This analysis encompassed 910 used PAS sites, 1,820 unused  
 1014 PAS sites, and  $\sim 5.7$  million genomic methylation sites.
- 1015 B. Cell-type specific Methylation State around IPA Sites in B cells. Metagene analysis of DNA  
 1016 methylation within  $\pm 2500$  bp of PAS in B cells revealed distinct epigenetic profiles between  
 1017 used (orange) and unused (navy) PAS loci. Used PAS sites exhibited higher methylation  
 1018 levels compared to unused PAS sites across the entire genomic window (Kolmogorov-

1019 Smirnov test,  $P - \text{value} \leq 5.14 \times 10^{-177}$  ). This analysis included 679 used PAS sites, 1,358  
1020 unused PAS sites, and ~3.4 million.

1021 C. Cell-Type Specific Methylation Patterns at IPA Sites in CD14<sup>+</sup> Cells. Metagene analysis of  
1022 DNA methylation within  $\pm 2500$  bp of PAS in CD14<sup>+</sup> cells revealed that used PAS sites  
1023 (orange) are consistently flanked by higher methylation levels compared to unused PAS  
1024 sites (navy) (Kolmogorov-Smirnov test,  $P - \text{value} \leq 3.94 \times 10^{-244}$ ). This analysis included  
1025 325 used PAS sites, 650 unused PAS sites, and ~2.95M methylation sites.

1026 D. DNA Methylation Landscapes at IPA Sites in AML. Metagene analysis of DNA methylation  
1027 within  $\pm 2500$  bp of PAS sites in AML demonstrates that used PAS sites (orange) are  
1028 consistently associated with higher levels of DNA methylation compared to unused PAS  
1029 sites (navy) (Kolmogorov-Smirnov test,  $P - \text{value} < 10^{-244}$ ). This analysis encompassed 552  
1030 used PAS sites, 1,104 unused PAS sites, and ~2.4 million methylation sites.

1031 E. Rolling Window Analysis of DNA methylation in IPA and Non-IPA introns. Mean DNA  
1032 methylation levels (y-axis) were calculated for IPA introns (violet), non-IPA introns within IPA  
1033 genes (yellow), and introns from genes without IPA events (brown) using a sliding window of  
1034 100 bp along the intron length (x-axis) with a 50 bp step size. This analysis included 696 IPA  
1035 introns, 1,758 non-IPA introns from IPA genes, and 1,303 introns from non-IPA genes,  
1036 encompassing ~5.7 million methylation sites in total. IPA introns consistently showed  
1037 significantly higher methylation levels across their entire length compared to both non-IPA  
1038 introns and introns from non-IPA genes (Kolmogorov-Smirnov test,  $P - \text{value} \leq 2 \times 10^{-16}$   
1039 for both comparisons).

1040 F. Heatmap of IPA site usage and mean methylation in the introns with differentially methylated  
1041 sites (GSE184314). The left heatmap shows the mean methylation in the introns coming  
1042 from differentially methylated sites and right heatmap shows the IPA site usage in the  
1043 corresponding introns. Each row represents a unique intron ( $n = 101$ ), while columns  
1044 correspond to individual patient samples ( $n = 10$ ). Rows are color-coded to indicate the

1045 differential methylation status in the introns: Hypermethylation (Green) (n = 38) and  
1046 Hypomethylation (Orange) (n = 63). Columns are color-coded to indicate the sample groups:  
1047 Control (Blue) (n = 5) and Diseased (Pink) (n = 5). Methylation is quantified on a 0-1 scale  
1048 (blue) with color intensity reflecting the degree of methylation and IPA site usage is  
1049 quantified on a 0-1 scale (red) relative to full-length isoform usage, with color intensity  
1050 reflecting the degree of IPA site usage.

1051 G. Density Plot of IPA site usage in Introns with Differentially Methylated Sites. Density plots  
1052 compare IPA site usage in introns harboring differentially methylated sites, separated into  
1053 hypermethylated (left panel) and hypomethylated (right panel) groups as identified in **Fig.**  
1054 **6F**. In the hypermethylated group (n = 380), IPA site usage differs significantly between  
1055 Control (blue) and Diseased (pink) samples (Wilcoxon paired test,  $P$  – value  $\leq$   
1056  $3 \times 10^{-10}$ ), whereas in the hypomethylated group (n = 630), no significant difference is  
1057 observed (Wilcoxon paired test,  $P$  – value  $< 0.56$ ).

1058 H. Heatmap of IPA site usage and mean methylation in the introns with differentially used IPA  
1059 sites (GSE184314). The left heatmap shows the IPA site usage in the introns coming from  
1060 differentially used IPA sites. and right heatmap shows the methylation in the corresponding  
1061 introns. Each row represents a unique intron (n = 102), while columns correspond to  
1062 individual patient samples (n = 10). Rows are color-coded to indicate the differential IPA site  
1063 usage status in the introns: highly used (dark green; n = 36) and lowly used (brown; n = 66).  
1064 Columns are color-coded to indicate the sample groups: Control (blue; n = 5) and Diseased  
1065 (pink; n = 5). Methylation is quantified on a 0-1 scale (blue) with color intensity reflecting the  
1066 degree of methylation and IPA site usage is quantified on a 0-1 scale (red) relative to full-  
1067 length isoform usage, with color intensity reflecting the degree of IPA site usage.

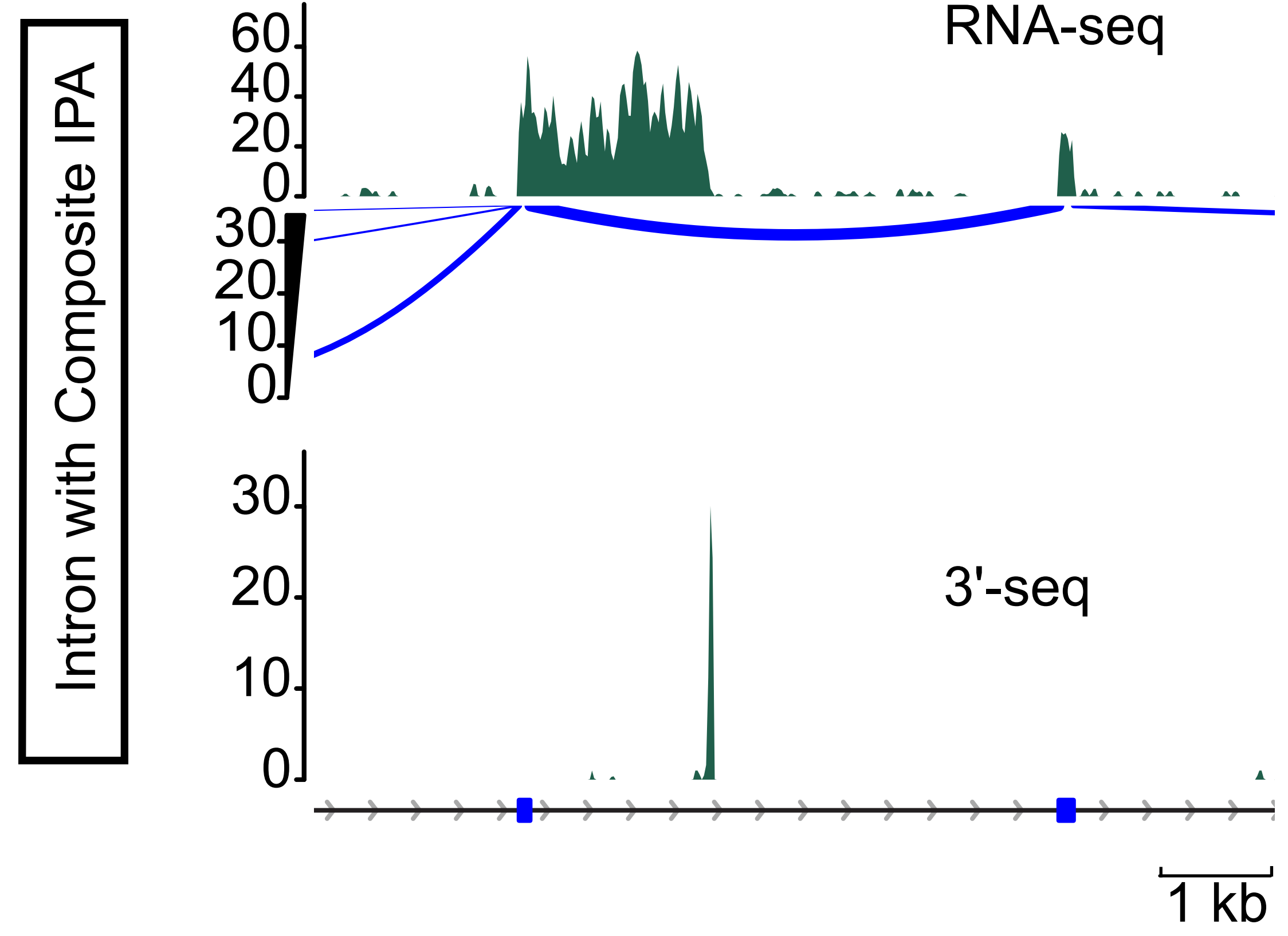
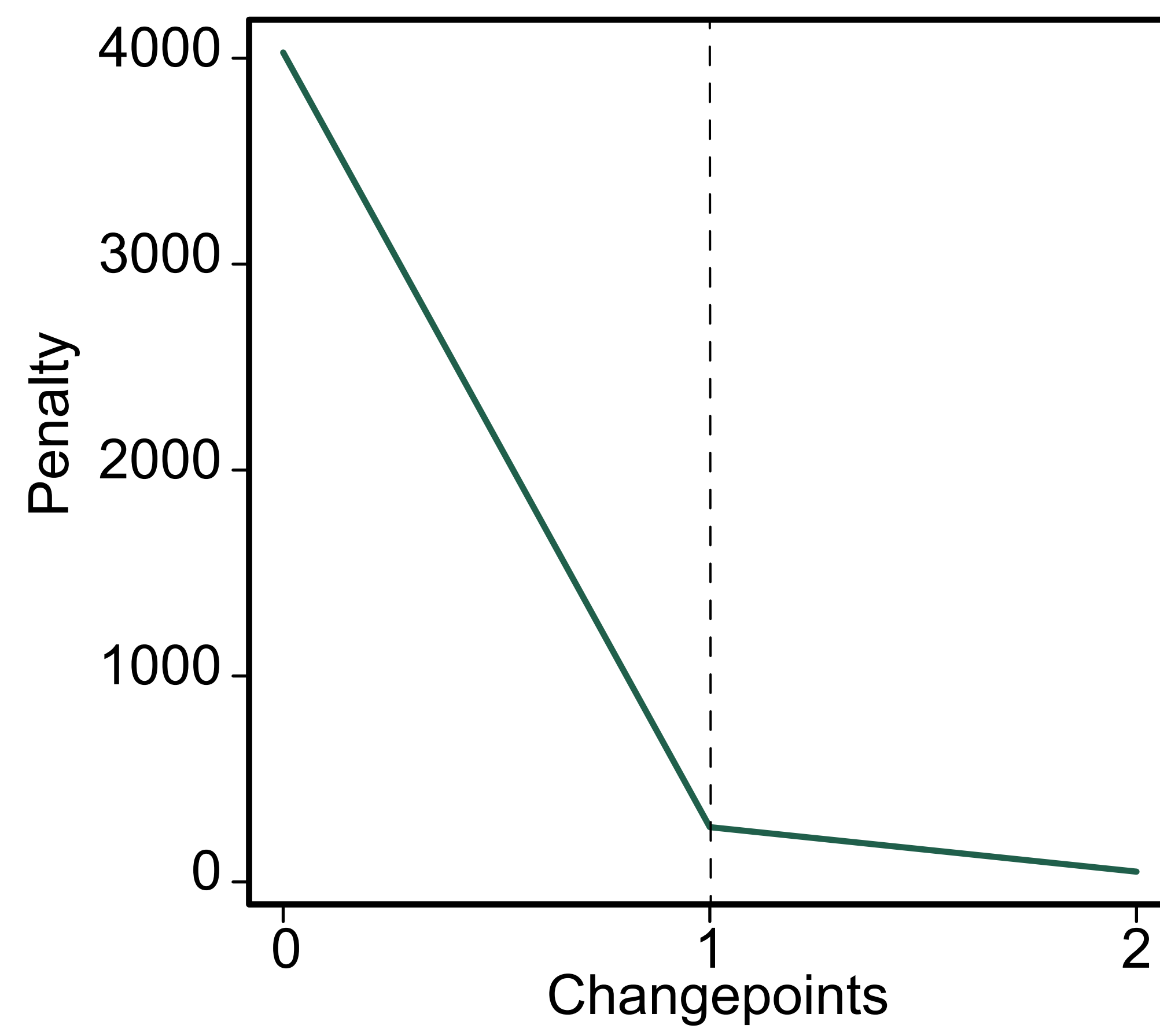
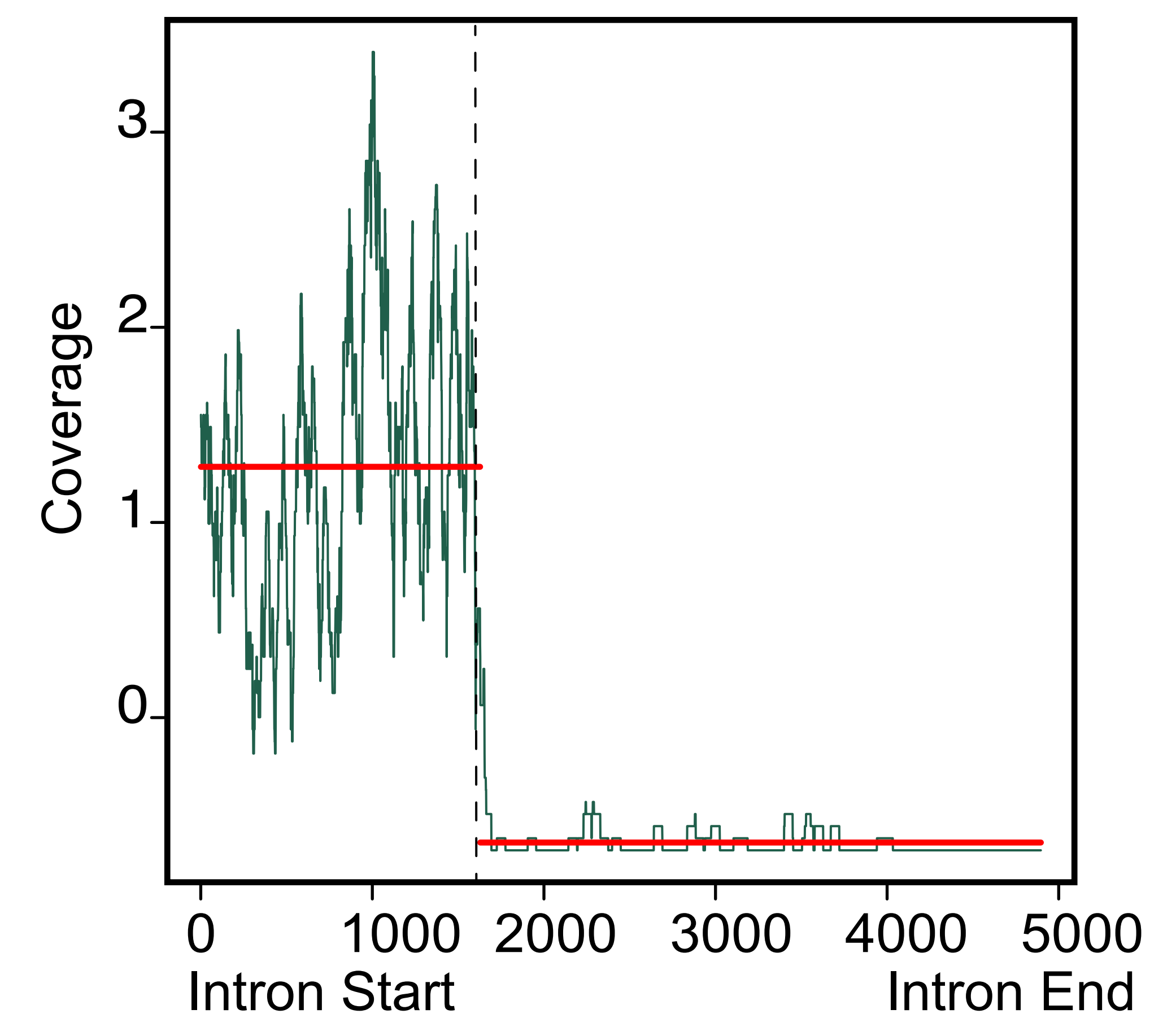
1068 I. Density Plot of Methylation Levels in Introns with Differentially Used IPA Sites. Density plots  
1069 illustrate methylation levels in introns containing differentially used IPA sites, divided into  
1070 lowly used (left panel) and highly used (right panel) groups as defined in **Fig. 6H**. For lowly

1071 used IPA sites ( $n = 660$ ), a modest but significant difference in methylation is observed  
1072 between control (blue) and diseased (pink) samples (Wilcoxon paired test,  $P$  – value <  
1073 0.15). In contrast, highly used IPA sites ( $n = 360$ ) show a pronounced and highly significant  
1074 methylation difference between groups (Wilcoxon paired test,  $P$  – value  $\leq 2.1 \times 10^{-10}$ ).

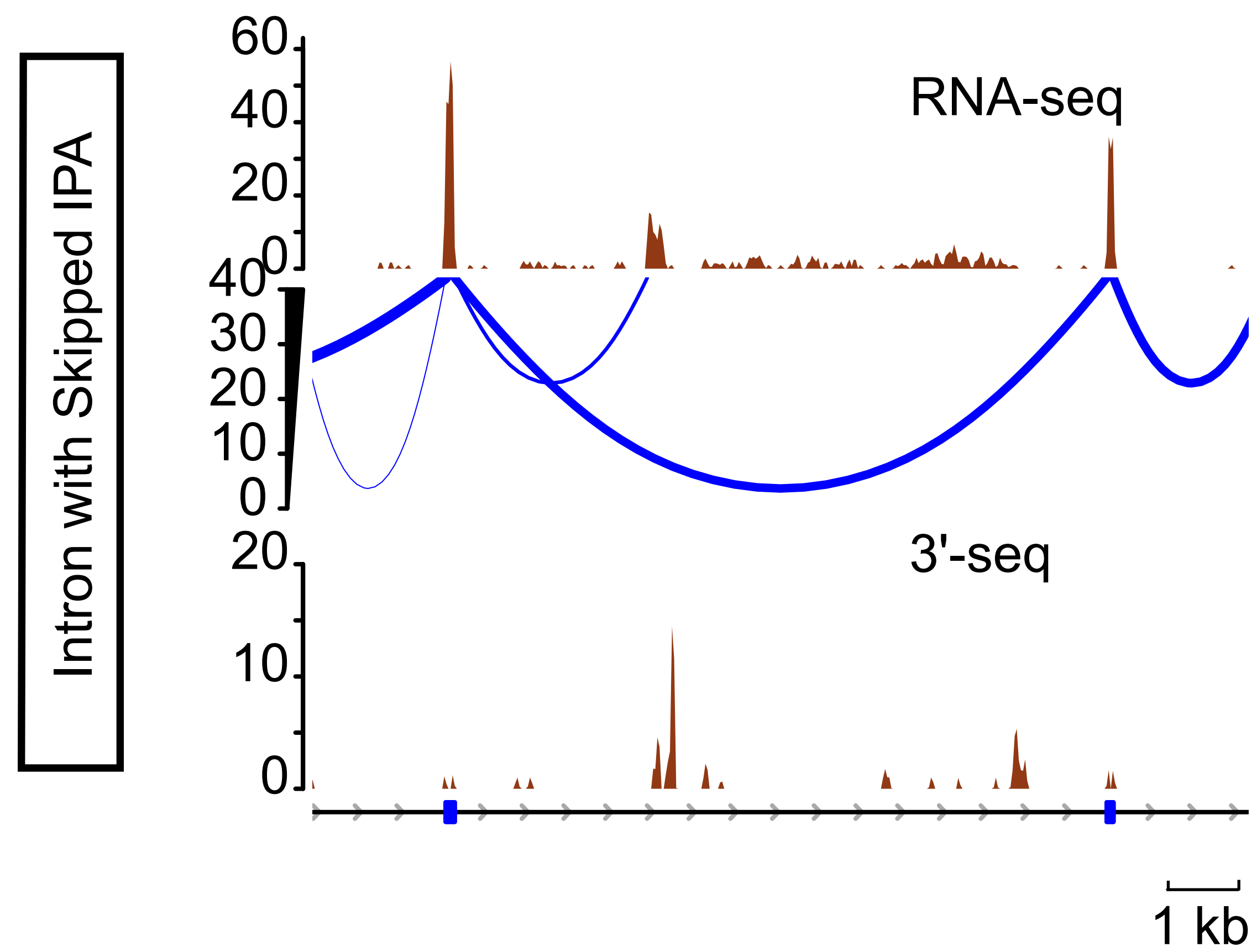
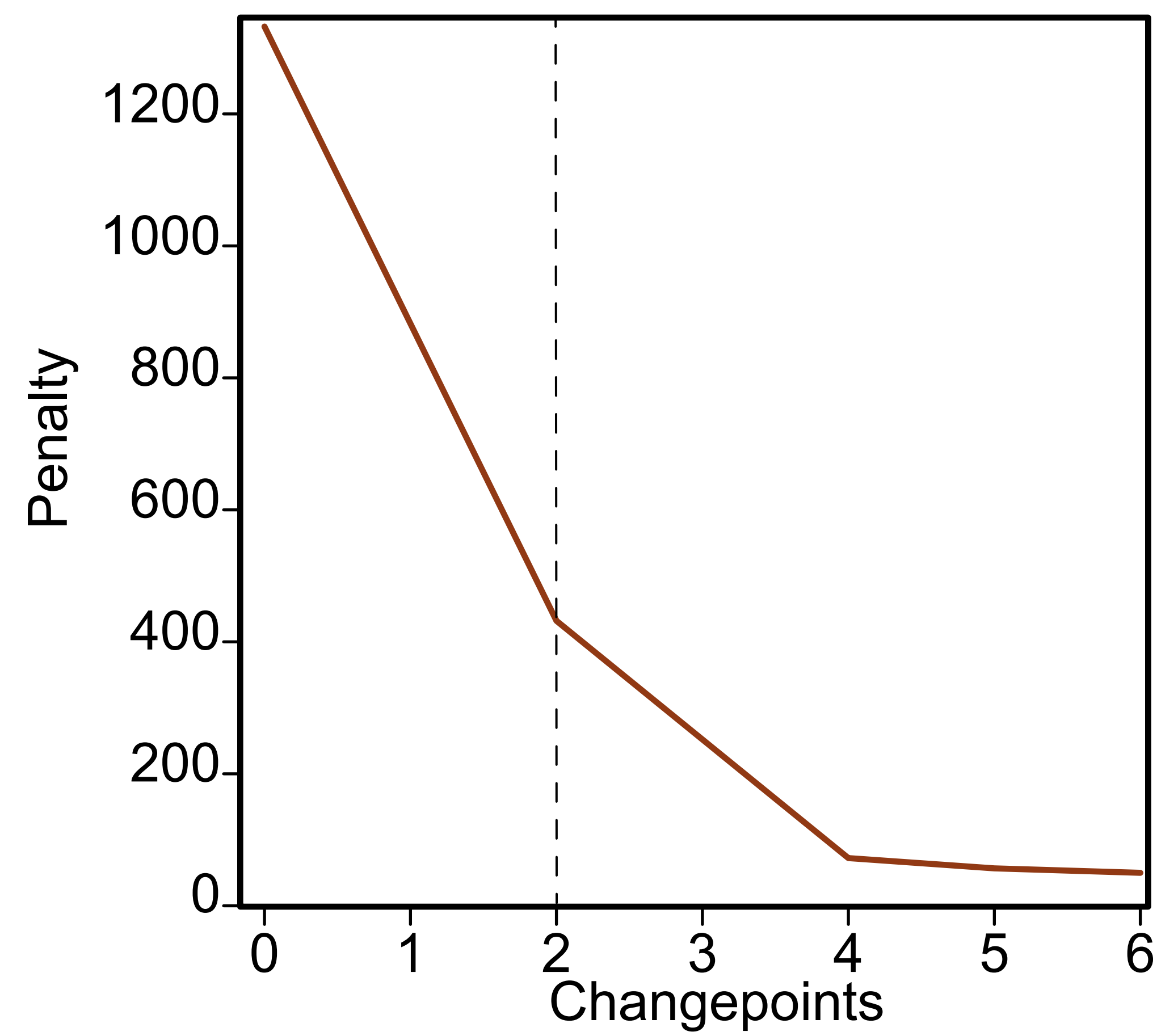
1075

**A**

KDM4B (Chr19: 5075471–5084366(+))

**B****C****D**

TPP2 (Chr13: 102602922–102616100(+))

**E****F**