

1 **A Survey on Methods for Predicting Polyadenylation Sites from**
2 **DNA Sequences, Bulk RNA-seq, and Single-cell RNA-seq**

3 Wenbin Ye¹, Qiwei Lian^{1,2}, Congting Ye³, Xiaohui Wu^{1,*}

4

5 ¹ *Pasteurien College, Soochow University, Suzhou 215000, China*

6 ² *Department of Automation, Xiamen University, Xiamen 361005, China*

7 ³ *Key Laboratory of the Ministry of Education for Coastal and Wetland Ecosystems,*
8 *College of the Environment and Ecology, Xiamen University, Xiamen 361005, China*

9 * Corresponding author.

10 E-mail: xhwu@suda.edu.cn (Wu X).

11 **Running title:** *Ye W et al / Survey on Poly(A) Site Prediction*

12

13 Word count: 9280

14 Keyword count: 5

15 Figure count: 3

16 Table count: 1

17 Supplementary Figures: 0

18 Supplementary Tables: 5

19 Supplementary Files: 1

20 Reference Count: 133

21 Article Title: 99 characters

22 Running Title: 29 characters

23 Abstract: 187 words

24

25

26 **Abstract**

27 Alternative polyadenylation (APA) plays important roles in modulating mRNA stability,
28 translation, and subcellular localization, and contributes extensively to shaping eukaryotic
29 transcriptome complexity and proteome diversity. Identification of poly(A) sites (pAs) on
30 a genome-wide scale is a critical step toward understanding the underlying mechanism of
31 APA-mediated gene regulation. A number of established computational tools have been
32 proposed to predict pAs from diverse genomic data. Here we provided an exhaustive
33 overview of computational approaches for predicting pAs from DNA sequences, bulk
34 RNA-seq data, and single-cell RNA-seq (scRNA-seq) data. Particularly, we examined
35 several representative tools using RNA-seq and scRNA-seq data from peripheral blood
36 mononuclear cells and put forward operable suggestions on how to assess the reliability of
37 pAs predicted by different tools. We also proposed practical guidelines on choosing
38 appropriate methods applicable to diverse scenarios. Moreover, we discussed in depth the
39 challenges in improving the performance of pA prediction and benchmarking different
40 methods. Additionally, we highlighted outstanding challenges and opportunities using new
41 machine learning and integrative multi-omics techniques and provided our perspective on
42 how computational methodologies might evolve in the future for non-3' UTR, tissue-
43 specific, cross-species, and single-cell pA prediction.

44 **Keywords:** Polyadenylation; Predictive modeling; RNA-seq; Single-cell RNA-seq;
45 Machine learning

46

47

48 **Introduction**

49 Precursor mRNA (pre-mRNA) polyadenylation is an essential two-step event in the post-
50 transcriptional regulation of gene expression, which involves the cleavage of the pre-
51 mRNA at the poly(A) site (pA) followed by the addition of an untemplated stretch of
52 adenosines [1, 2]. The selective use of pAs of a single gene, termed alternative
53 polyadenylation (APA), can generate a diversity of isoforms with different 3' ends and/or
54 encode distinct proteins [3, 4]. APA plays important roles in modulating mRNA stability,
55 translation, and subcellular localization, which contributes extensively to shaping
56 eukaryotic transcriptome complexity and proteome diversity. APA is a widespread
57 regulatory mechanism in eukaryotes, which has been observed in more than 70% of
58 mammalian and plant genes [5-11]. APA is highly tissue specific and dynamically
59 modulated in various conditions, cell types, and/or states [2, 12]. Specific APA programs
60 have been implicated in diverse biological processes and diseases, such as cell activation,
61 proliferation, neurodegenerative disorders, and cancer [3, 4, 13-20]. Given the functional
62 significance of APA, identification and/or quantification of pAs on a genome-wide scale is
63 crucial and may be the first step in understanding the underlying mechanism of APA-
64 mediated gene regulation.

65 Early studies, dating back to the 1990s, predict pAs using conventional machine
66 learning (ML) models like support vector machine (SVM) [21-25], which distinguish
67 whether a nucleotide sequence contains a pA using a variety of hand-crafted features
68 (**Figure 1A**). In recent years, deep learning (DL) models [26-29] have been shown to
69 provide better performance than traditional ML methods, owing to their great ability for
70 direct and automatic feature extraction and high scalability with large amount of genomic
71 data (Figure 1B). With the advance of next generation sequencing (NGS) technologies,
72 experimental protocols have been designed to capture 3' ends of mRNAs for direct
73 profiling of genome-wide pAs (Figure 1C), such as DRS [10, 30], 3P-Seq [7, 31], 3'READs
74 [11], PAT-seq [32], TAIL-seq [33, 34], and several others (reviewed in [35-37]). Although

75 these 3' end sequencing (3' seq) approaches are powerful and highly sensitive in detecting
76 the precise locations of pAs, even for lowly expressed genes, they are too technically
77 demanding and costly to be widely applied in genomic research. Alternatively, a myriad of
78 computational tools [38-41] have been developed for identifying and quantifying pAs by
79 leveraging the explosively growing RNA sequencing (RNA-seq) data from diverse
80 biological conditions, cell types, individuals, and organisms (Figure 1D). In recent years,
81 the single-cell RNA-seq (scRNA-seq) techniques, particularly those 3' tag-based protocols
82 such as CEL-seq [42] and 10x Chromium [43], provide great potential to explore dynamics
83 of APA usage during the process of cellular differentiation. Accordingly, a wide spectrum
84 of tools have been proposed to profile APA from diverse scRNA-seq datasets at cell-type
85 or even single-cell resolution [44-46] (Figure 1E).

86 The tsunami of genomic data especially bulk and single-cell RNA-seq data and the
87 emergence of ensemble deep learning methodologies have revolutionized computational
88 methods for detecting pAs from diverse kinds of data. In the past decade, a few literature
89 reviews have involved the computational tools for bioinformatic analysis of APA. In 2015,
90 our group summarized computations tools for predicting pAs from DNA sequences and 3'
91 seq methods for mapping pAs [37]. Szkop and Nobel [47] described experimental methods
92 for probing 5' UTRs and 3' UTRs, and listed computational methods for discovering
93 alternative transcription start sites (TSSs) and pAs from microarray and RNA-seq. Yeh et
94 al. [48] reviewed experimental methods and technologies for studying APA, and briefly
95 listed seven RNA-seq tools for analyzing APA dynamics in tabular form. Chen et al. [49]
96 comprehensively reviewed 3' seq methods for probing pAs, while their review did not
97 cover the computational tools for APA analysis. Gruber and Zavolan [12] highlighted the
98 importance of APA in health and disease, and briefly listed computational resources for
99 studying APA in a table, including four pA databases, two databases of RBP binding motifs,
100 eight RNA-seq tools for identifying and/or quantifying pAs, and three tools for APA
101 analysis. Our group [50] benchmarked 11 tools for predicting pAs or dynamic APA events

102 from RNA-seq data. Another benchmark study [51] benchmarked five tools for RNA-seq
103 and compared their performance with 3' seq, Iso-Seq, and PacBio single-molecule full-
104 length RNA-seq method. Ye et al. [52] briefly summarized three computational methods
105 for detecting APA dynamics from diverse single cell types. Zhang et al. [53] focused on the
106 APA regulation in cancer, and briefly listed 14 computational tools for detecting APA.
107 Kandhari et al. [54] highlighted the emerging role of APA as cancer biomarkers and
108 provided an overview of existing relevant experimental and computational methods.
109 However, these two reviews [53, 54] did not distinguish among the prediction of pAs,
110 detection of APA dynamics, and analysis of APA. For example, APALyzer [55] and
111 movAPA [56] listed in these reviews are actually toolkits for analyzing APA rather than
112 detecting APA dynamics or pAs, which are different from other tools they listed such as
113 DaPars [39] or APATrap [40]. Generally, although the above reviews have provided detailed
114 overviews of the progress in the complex yet fruitful APA field, none of them has
115 exhaustively summarized available tools for different kinds of data in this field, particularly
116 the emerging DL-based methods and methods for scRNA-seq. Moreover, most reviews
117 only briefly listed tools without delicate summary and sorting, which makes it difficult for
118 the scientific community to decide desirable method for their data analysis. In this review,
119 we described the principles of identifying pAs from different kinds of data and provide an
120 extensive overview of available computational approaches. We catalogued these methods
121 into different categories in terms of the underlying principles of the predictive models and
122 the data they used, and summarized their performance and characteristics such as
123 algorithms, features, and data used in the predictive model. Particularly, we examined
124 several representative tools using RNA-seq and scRNA-seq data from peripheral blood
125 mononuclear cells and put forward operable suggestions on how to assess the reliability of
126 pAs predicted by different tools. We also describe several notes on how to conduct
127 objective benchmark analysis for these massive number of tools. Moreover, we propose
128 practical recommendations on choosing appropriate methods for different scenarios and

129 discussed implications and future directions. Additionally, we highlight outstanding
130 challenges and opportunities using new machine learning and integrative multi-omics
131 techniques. Lastly, we provide our perspective on how computational methodologies might
132 evolve in the future for pA prediction, including non-3' UTR, tissue-specific, cross-species,
133 and single-cell pA prediction.

134 **Computational approaches for pA prediction**

135 **Methods for predicting pAs from DNA sequences**

136 The key trigger for cleavage and polyadenylation is the set of *cis*-regulatory elements
137 surrounding a pA, including A[A/U]UAAA hexamer or variant thereof, the UGUA element,
138 upstream and downstream U-rich elements, and downstream GU-rich elements [57]. Since
139 poly(A) signals, the core AAUAAA and its variants, are in the vicinity of most mammalian
140 pAs, the identification of the poly(A) signal (PAS) is usually regarded as an alternative to
141 determine the potential position of a pA. In this review, we refer to the task of predicting
142 pAs or PASs as the "pA identification problem". During the past few decades, a wide range
143 of computational approaches have been proposed to predict pAs from DNA sequences
144 using experimental and *in silico* mapping of 3'-end expressed sequence tags (ESTs) (Files
145 S1 and S2).

146 *Methods based on traditional machine learning models*

147 Earlier studies established traditional ML models to classify a sequence as containing a pA
148 or not, using various algorithms such as discriminant functions [21, 22, 58], hidden Markov
149 model (HMM) [23], SVM [24, 59], Bayesian network [60], artificial neural network and
150 random forests [61], and combined classifiers [25, 62] (**Figure 2** and File S2). The machine
151 learning frameworks of these methods are similar, except that different classification
152 models were employed and/or diverse hand-crafted sequence features were compiled (File
153 S1). As ML models rely heavily on manually designed features and the poly(A) signal of
154 human/animal is considerably different from that of other species like plants or
155 *Saccharomyces cerevisiae* (yeast) [37, 63], these ML-based methods can be divided into

156 two categories according to the applicable species (File S1): i) methods that are applicable
157 to human or animals, including POLYAH [21], Polyadq [58], ERPIN [23], Poly(A) Signal
158 Miner [64], Polya_svm [24], PolyApred [59], POLYAR [22], Chang's model [65], Dragon
159 PolyA Spotter [61], Xie's model [66], and Omni-PolyA [25]; ii) methods that are applicable
160 to other species, including the Graber's method [67] for yeast, POLYA [68] for
161 *Caenorhabditis elegans*, PASS [69, 70], PAC [60] and PASPA [71] for plants, and Wu's
162 model for *Chlamydomonas Reinhardtii* [62]. These methods utilize diverse sequence
163 features around pAs for pA prediction (File S1). The most commonly used features are
164 position weight matrix for the poly(A) motifs, distance between motifs, and k-gram
165 nucleotide acid patterns [21, 23, 24, 58, 59]. With the increase of the prior knowledge of
166 DNA sequences, more carefully hand-crafted features were derived, such as Z-curve [60],
167 RNA secondary structures [62, 65], physico-chemical, thermodynamic and statistical
168 characteristics [61], the term frequency-inverse document frequency weight [62], and
169 spectral latent features extracted by HMM [66]. Particularly, since the significance of
170 poly(A) signal is different in pAs with different strengths, a few studies divided pAs into
171 sub-groups based on the expression level [22] or pattern assembly [62], and then predicted
172 pAs in each group. In terms of the availability and ease of use of tools, several tools were
173 presented as website (Figure 2), which is particularly convenient for users with little
174 program skill. However, since these tools were generally developed many years ago, the
175 programming languages of many tools are outdated, such as Fortran or Perl, and many tools
176 are no longer available or maintained.

177 *Methods based on deep learning models*

178 Despite considerable progress has been made, the overall accuracy and generalizability of
179 traditional ML-based methods remain moderate due to the limited experimentally verified
180 pAs in the early years and the lack of prior domain knowledge to finely design and acquire
181 useful features. In recent years, DL-based methods are emerging rapidly (File S2 and
182 Figure 2), which directly learn hidden features from input nucleotide sequences in a data-

183 driven manner, without knowing any prior knowledge of sequence features. Most methods
184 use convolution neural networks (CNNs), including deepPolyA [72], Conv-Net [73],
185 DeeReCT-PolyA [26], DeepPASTA [28], DeepGSR [27], and APARENT [29]. Other deep
186 learning techniques were also utilized, such as the recurrent neural network (RNN)
187 employed in DeepPASTA [28], a hybrid model with four logistic regression models and
188 eight neural networks used in HybPAS [74], and self-attention mechanisms used in
189 SANPolyA [75] and PASNet [76]. All of these tools were implemented using DL
190 frameworks in Python. In addition to pA prediction, several methods can be utilized for
191 multiple tasks. For example, Conv-Net [73] is capable of inferring pA selection and
192 predicting pathogenicity of polyadenylation variants. DeepPASTA [28] can be used for the
193 prediction of the most dominant pA of a gene in a given tissue and the relative dominance
194 of APA sites in a gene. DeepGSR [27] is able to predict genome-wide and cross-organism
195 genomic signals such as translation initiation sites. APARENT [29] can also be utilized for
196 the quantification of the impact of genetic variants on APA. Different from hand-picked
197 features used in ML-based methods, one-hot encoding features without needing fine
198 feature engineering are widely used in DL-based methods, however, DL-based models are
199 generally of poor interpretability. To enhance the interpretability, several methods provide
200 additional function for visualization of signals. Xia et al. [26] showed the interpretability
201 of their DeeReCT-PolyA model by transforming convolutional filters into sequence logos
202 for the comparison between human and mouse. In APARENT [29], features learned across
203 all network layers were visualized, which can reveal *cis*-regulatory elements known to
204 recruit APA regulators and new sequence determinants of polyadenylation. In addition to
205 performance improvement, DL-based methods have two significant advantages over ML-
206 based methods, the higher generalizability for different species and the higher scalability
207 with large amount of data. For example, DeeReCT-PolyA [26] is an interpretable and
208 transferrable CNN model for recognition of 12 PAS variants, which enables transfer
209 learning across datasets and species. APARENT [29] was trained using isoform expression

210 data from more than three million synthetic APA reporters.

211 **Methods for predicting pAs from bulk RNA-seq data**

212 Methods that predict pAs only from DNA sequences conspicuously fail to consider *in vivo*
213 expression. RNA-seq has become an indispensable approach for transcriptome profiling in
214 diverse biological samples and a number of methods have been proposed for identifying
215 sample-specific pAs from RNA-seq (File S3). Our group previously benchmarked 11
216 representative methods for predicting pAs and/or dynamic APA events from RNA-seq [50].
217 Here we focus on prediction of pAs rather than dynamic APA events. We collected relevant
218 methods summarized in our previous review [50] as well as newly emerging methods, and
219 divided these methods into five categories according to their underlying strategies.

220 *Methods that interrogate non-templated poly(A)-capped reads*

221 RNA-seq data contain a small fraction (~0.1%) of non-templated poly(A) tail-containing
222 reads (hereinafter referred to as poly(A) reads) [47], which can be considered as direct
223 evidence for polyadenylation. By interrogating poly(A) reads, an early study [77] identified
224 ~8000 novel pAs in *Drosophila melanogaster* from a total of 1.2 billion RNA-seq reads.
225 Several other methods, such as KLEAT [78] and ContextMap 2 [79], not only employed
226 direct evidence from poly(A) reads but also incorporated transcript assembly to identify
227 pAs. These poly(A) read-based approaches have the advantage to determine the precise
228 locations of pAs, however, it is still challenging to discover pAs of weakly expressed
229 transcripts due to the decreased read coverage near the 3' end and the low yield of poly(A)
230 reads.

231 *Methods based on transcript assembly*

232 Another series of approaches identify pAs from inferred alternative 3' UTRs by compiling
233 transcript structures from RNA-seq, including PASA [80], Scripture [81], 3USS [82], and
234 ExUTR [83]. These transcriptome assembly-assisted methods deduce gene models first
235 using transcriptome assembly tools, and then identify 3' UTRs that are absent in the
236 deduced gene models, which rely heavily on assembled gene structures. It is widely

237 accepted that transcriptome assembly from RNA-seq is a rather difficult and
238 computationally demanding task, and it is more challenging to precisely determine 3' UTRs,
239 especially for lowly expressed genes, due to 3' biases of read coverage inherent in RNA-
240 seq. Therefore, the performance of these methods is inevitably hindered by potential
241 limitations of existing transcriptome assembly tools.

242 *Methods that rely on prior annotations of pAs*

243 During the last decade, numerous experimental techniques have been developed to direct
244 sequence 3' ends of mRNAs, such as 3' T-fill [84], 3'READs [11], TAIL-seq [33, 34], to
245 name a few (Figure 1C). Accordingly, several pA databases built upon 3' seq data of diverse
246 species were continuously released, including PolyA_DB 3 [85], PolyAsite 2.0 [8], and
247 PlantAPAdb [86]. These databases provide a large number high-confidence pAs, which can
248 be used for establishing pA prediction models and evaluating pA prediction results. It is
249 thus naturally to incorporate annotated pAs for predicting pAs from RNA-seq. Several
250 methods, including QAPA [38], PAQR [87], and APA-scan [88], that rely on pre-defined
251 pA annotations were proposed for predicting pAs from RNA-seq. For these methods, the
252 quality of annotated pAs is particularly critical. Most studies establish a comprehensive
253 compendium of well-annotated pAs by merging non-redundant annotations from diverse
254 sources. By combining priori annotated pAs with RNA-seq, the quality of predicted pAs
255 can be greatly improved. However, currently available pA databases are far from complete
256 and limited to only a few well-studied species, such as human, mouse, and *Arabidopsis*
257 *thaliana*, consequently, these tools are not capable of detecting novel pAs beyond existing
258 poly(A) annotations.

259 *Methods that infer pAs by detecting significant changes in RNA-seq read density*

260 Majority of recent approaches predict pAs by modelling read density changes in terminal
261 exons, including GETUTR [89], IsoSCM [90], DaPars/DaPars2 [39, 91, 92],
262 EBChangePoint [93], APAtrap [40], TAPAS [41], moutainClimber [94], and IPAFinder
263 [95]. According to our previous benchmark on 11 tools for RNA-seq [50], TAPAS generally

264 obtained higher sensitivity than other tools across different datasets. Of note, unlike most
265 methods that require at least two samples for change point detection, mountainClimber [94]
266 is a *de novo* cumulative-sum-based approach, which runs on a single RNA-seq sample and
267 simultaneously recognizes multiple TSSs or APA sites in a transcript. Using
268 mountainClimber, Cass and Xiao analyzed 2,342 GTEx samples from 36 tissues of 215
269 individuals and found 75% of genes exhibited differential APA across tissues [94].
270 Different from most pA prediction tools focusing mainly on 3' UTR, IPAFinder was
271 specifically proposed for identifying intronic pAs from RNA-seq [95]. Zhao et al. applied
272 IPAFinder to pan-cancer datasets across six tumor types and discovered 490 recurrent
273 dynamically changed intronic pAs [95]. Methods falling within this category rely on the
274 detection of read density fluctuations which require sufficient read coverage in terminal
275 exons to detect APA sites. It is worth noting that data pre-processing (normalization or
276 smoothing) is particularly important for reducing technical biases caused by non-biological
277 variability [47]. Particularly, some methods, such as APAtrap and DaPars, re-define
278 terminal exon boundaries based on RNA-seq read coverage before identifying pAs, which
279 are capable of detecting pAs in previously unannotated regions.

280 *Methods based on machine learning models*

281 In recent years, some newly emerging methods employ traditional ML or DL model to
282 identify pAs from RNA-seq, including TECtools [96], IntMAP [97], Terminitor [98], and
283 Aptardi [99]. TECtools [96] first identifies terminal exons and transcript isoforms ending
284 at known intronic pAs. Then a model was trained based on the aligned RNA-seq data for
285 distinguishing terminal exons from internal exons and background regions, using diverse
286 features reflecting differences in read coverage of these regions. TECtool can also be
287 applied on scRNA-seq, which first pools reads of all cells to infer new transcripts and then
288 quantify each transcript in individual cells. IntMAP [97] leverages one unified ML
289 framework to combine the information from RNA-seq and 3' seq to quantify different 3'
290 UTR isoforms using a global optimization strategy. Terminitor [98] is based on a deep

291 neural network for three-label classification problem, which can determine whether an
292 input sequence contains a pA with poly(A) signal, a site without poly(A) signal, or non-pA.
293 Aptardi [99] is a multi-omics approach based on bidirectional long short-term memory
294 recurrent neural network (biLSTM), which predicts pAs by leveraging DNA sequences,
295 RNA-seq, and the predilection of transcriptome assemblers.

296 **Methods for predicting pAs from single-cell RNA-seq**

297 Single-cell RNA-seq is a powerful high-throughput technique for interrogating
298 transcriptome of individual cells and measuring cell-to-cell variability in transcription
299 [100]. Particularly, several 3' tag-based scRNA-seq methods enriching for mRNA 3' ends
300 via poly(A) priming, such as CEL-seq [42], Drop-seq [101], and 10x Chromium [43],
301 provide great potential to dissect APA at single-cell resolution. However, the extremely
302 high dropout rate and cell-to-cell variability inherent in scRNA-seq makes it difficult to
303 directly apply bulk RNA-seq methods to scRNA-seq data. During the last few years, a wide
304 range of computational approaches specifically designed for pA identification from
305 scRNA-seq have emerged (File S4 and Figure 2). We divided these methods into three
306 categories according to their underlying strategies.

307 *Methods based on peak calling*

308 The peak calling strategy is widely used by most methods for pA identification from
309 scRNA-seq, including scAPA [102], polyApipe
310 (<https://github.com/MonashBioinformaticsPlatform/polyApipe>), Sierra [44], scAPATrap
311 [45], SAPAS [103], and SCAPE [104]. The underlying principle of these methods is that
312 aligned reads from 3' tag-based scRNA-seq accumulate to form peaks at genomic intervals
313 upstream of pAs [102]. In scAPA [102], a set of non-overlapping 3' UTRs is first defined
314 from the genome annotation and then peaks within 3' UTRs are identified using an existing
315 peak calling tool. As adjacent pAs may situate in a single peak, the Gaussian finite mixture
316 model was implemented in scAPA to split large peaks into smaller ones. polyApipe is a
317 pipeline for identifying pAs from 10x Chromium scRNA-seq, which defines peaks of

318 polyA-containing reads. Sierra [44] employed the splice-aware peak calling based on
319 Gaussian curve fitting to determine potential peaks with pAs and then the peaks were
320 annotated and quantified in individual cells. Our group proposed scAPATrap [45] for
321 identifying and quantifying pAs in individual cells from 3' tag-based scRNA-seq.
322 scAPATrap incorporates a genome-wide sensitive peak calling strategy and poly(A) read
323 anchoring, which can accurate locate pAs without using prior genome annotation, even for
324 those with very low read coverage. Yang et al. proposed SAPAS for identifying pAs from
325 poly(A)-containing reads and quantifying pAs in peak regions determined by a parametric
326 clustering algorithm [103]. They further applied SAPAS to the scRNA-seq data of
327 GABAergic neurons and detected cell type-specific APA events and cell-to-cell modality
328 of APA for different GABAergic neuron types. Very recently, Zhou et al. proposed the
329 SCAPE method based on a probabilistic mixture model for identification and quantification
330 of pAs in single cells by utilizing insert size information [104]. The parametric modeling
331 of peaks in most tools based on peak calling such as scAPA or Sierra may cause biases and
332 reduce statistical power in detecting APA events. Alternatively, ReadZS [105], an
333 annotation-free statistical approach, was proposed to characterize read distributions that
334 bypasses parametric peak calling and identify differential APA usages at single-cell
335 resolution among ≥ 2 cell types. ReadZS can not only detect pAs in normal peak shape,
336 but also identify distributional shifts that are not.

337 *Methods that rely on prior annotations of pAs*

338 In contrast to the peak calling-based methods used for *de novo* pA identification, a few
339 approaches identify pAs base on prior pA annotations, including MAPPER [106],
340 SCAPTURE [107], and scUTRquant [108]. Li et al. developed MAPPER [106] for
341 predicting pAs from both bulk RNA-seq and scRNA-seq data, which incorporates
342 annotated pAs in PolyA_DB 3 [85] and pools single cells of the same type to mimic
343 pseudo-bulk samples. MAAPER also provides a likelihood-based statistical framework for
344 analyzing APA changes and can identify common and distinct APA events in cell groups

345 from different individuals. The group of MAPPER later developed SCAPTURE [107]
346 which embedded a DL model DeepPASS for evaluating called peaks from scRNA-seq. The
347 DL model was trained by sequences shifting, using annotated pAs from PolyA_DB 3,
348 PolyA-seq, PolyASite 2.0 and GENCODE v39. The authors used SCAPTURE to profile
349 APA dynamics between COVID-19 patients and healthy individuals, and found the
350 preference of proximal pA usage in numerous immune response-associated genes upon
351 SARS-CoV-2 infection. Fansler et al. developed scUTRquant [108] for measuring 3' UTR
352 isoform expression from scRNA-seq, which relies on a cleavage site atlas established from
353 GENCODE annotation and a mouse Microwell-seq dataset of 400,000 single cells [109].

354 *Other methods for predicting pAs from scRNA-seq*

355 Additionally, some other methods do not use the peak calling strategy, including APA-Seq
356 [110] and scDaPars [46]. Levin et al. [110] designed the APA-seq approach to detect and
357 quantify pAs from CEL-seq, which interrogates the gene identity and poly(A) information
358 in the paired Read 1 and Read 2. Although APA-Seq is in principle applicable to other 3'
359 tag-based scRNA-seq methods, it may not be universally applied in practice in that only
360 sample barcodes rather than the whole 3' end sequence of the transcript are retained in Read
361 1 of many public scRNA-seq data [45]. Unlike most tools that are only applicable to 3' tag-
362 based scRNA-seq, scDaPars [46] that was proposed by the group of DaPars [39] can
363 identify and quantify APA events from either 3' tag (e.g., 10x Chromium) or full-length
364 (e.g., Smart-seq2) scRNA-seq. In the scDaPars pipeline, DaPars, a tool for identifying APA
365 events from bulk RNA-seq, was first adopted to calculate raw relative APA usage in
366 individual cells, and then a regression model was utilized to impute missing values in the
367 sparse single-cell APA usage matrix. By applying scDaPars to cancer and human endoderm
368 differentiation data, Gao et al. revealed cell type-specific APA regulation and detected
369 novel cell subpopulations that were not found in conventional gene expression analysis.

370 **Methods for APA analysis rather than pA prediction**

371 In addition to the task of pA prediction (hereinafter termed task 1), there are additional

372 tasks related to the bioinformatic analysis of APA, mainly including the prediction of
373 tissue-specific pAs (task 2), prediction of dominant pAs (task 3), prediction of APA site
374 switching (task 4), and other kinds of APA analysis (task 5). Although most tools described
375 in this review are developed for task 1, several tools are capable of performing multiple
376 tasks. For example, DeepPASTA [28] is able to perform tasks 1-3; Conv-Net [73] can
377 perform tasks 1/3. In this review, we focus only on tools that are applicable to task 1. Of
378 note, NGS-based techniques specially designed for probing pAs, generally known as 3' seq,
379 such as DRS [10, 30], 3P-Seq [7, 31], and 3'READs [11], are experimental methods rather
380 than computational methods for identifying pAs. Genome-wide pAs generated from 3' seq
381 are highly confident and are usually regarded as the true reference (i.e., prior information)
382 for building models or evaluating computational methods. These 3' seq methods are beyond
383 the scope of this review, while have been reviewed in several other reviews [12, 47, 49,
384 54]. In addition, we have briefly summarized tools or resources designed for APA analysis
385 rather than pA prediction in File S5. Tools such as DeeReCT-APA [111], polyA code [112],
386 and TSAPA [113] are not targeted at task 1 but for other tasks 2/3, such as predicting tissue-
387 specific pAs. Among the five tasks, detection of APA site switching (Task 4) is usually a
388 routine step involved in the analysis of RNA-seq or scRNA-seq. APA site switching reflects
389 the differential usage of APA sites between samples, which does not necessarily need the
390 prediction of pAs (task 1) as a prerequisite. Of note, there are other commonly used phrases
391 similar to 'APA site switching' mentioned in this review, such as differential APA site usage
392 [8, 39], 3' UTR shortening/lengthening [45, 102], and APA dynamics [39, 45, 99, 114].
393 Some approaches for RNA-seq, such as PHMM [115], ChangePoint [116], MISO [117],
394 and roar [118], directly discover APA site switching by detecting sudden change of read
395 density at terminal exons without identifying APA sites. Recently, several tools were
396 developed for scRNA-seq, such as SCUREL [119], scMAPA [120], and scDAPA [121].
397 For example, our group developed scDAPA [121] for characterizing differential usages of
398 APA in different cell types using 10x Chromium data, and found APA plays important role

399 in acute myeloid leukemia [114]. Additionally, some toolkits were developed for routine
400 analyses of APA (e.g., annotation and visualization, task 5) using annotated pAs and/or
401 RNA-seq, such as APAlyzer [55] and movAPA [56], while they are not capable of
402 predicting pAs. These diverse tools provide a wide range of complementary resources and
403 opportunities to address the more complex but fruitful field of APA.

404 **Discussion**

405 **Performance of pA prediction models**

406 At present, there are only a few benchmark studies that systematically evaluate the
407 performance of different tools. Previously, our group benchmarked 11 tools for RNA-seq
408 [50] and found that the sensitivity of some methods varied greatly among different species.
409 For instance, QAPA [38] performs the second best on human data, while it performs the
410 worst on mouse data. APAtrap [40] is the top performer for *Arabidopsis* data, while TAPAS
411 [41] performs the best on human or mouse data. Recently, Shah et al. [51] benchmarked
412 five tools for RNA-seq against 3' seq, Iso-Seq, and a full-length RNA-seq method and
413 found that pAs from 3' seq and Iso-Seq are more reliable than pAs predicted from RNA-
414 seq. They suggested that incorporating the RNA-seq prediction tool QAPA [38] with pA
415 annotations derived from 3' seq or Iso-Seq can reliably quantify APA dynamics across
416 conditions.

417 The performance of different tools described in the respective studies was summarized
418 in Files S1 to S4. Generally, for predicting pAs from DNA sequences, DL-based models
419 significantly outperformed ML-based methods and are more suitable for large-scale
420 analysis, owing to the good ability of automatic feature extraction and scalability for big
421 data analysis (Table S2). For example, DeepPASTA [28] has an area under the curve score
422 over 93% in predicting pAs on a DNA sequence dataset, which performed much better than
423 ML-based tools like PolyAR [22] or Dragon PolyA Spotter [61]. APARENT [29], based
424 on deep neural network, was trained on over three million synthetic APA reporter genes,
425 which overcomes inherent size limitations of traditional biological datasets. In contrast,

426 traditional ML-based methods like POLYAR [22] and Omni-PolyA [25] require a
427 considerable amount of prior knowledge and are unable to cope with the rapidly growing
428 data. In terms of the model generalizability, methods for RNA-seq or scRNA-seq are
429 generally applicable to different species if the reference genome and the genome annotation
430 are available. In contrast, the cross-species applicability of methods for DNA sequences is
431 more complex. Models applied to human are normally applicable to other mammals like
432 mouse due to similar poly(A) signals among mammals [57]. However, although most
433 models can be in principle trained using data from a different species, users need collect
434 training data from the other species which are not always available, and most models use
435 hand-crafted features that may not be generalized well across species. Recent techniques
436 like deep learning and transfer learning greatly enhance the generalizability of models.
437 Cross-species experiments have been performed for evaluating the generalizability of some
438 tools, such as DeepGSR [27] and Poly(A)-DG [122], and for these tools single model
439 trained over one species can be generalized well to datasets of other species without
440 retraining. We need to point out that, the evaluation results in a single study may be biased
441 and should be treated with caution, because different datasets and performance indicators
442 were used for the performance evaluation in different studies (Files S1 to S4). In the
443 following section of Conclusions and prospects, we also put forward several notes on how
444 to conduct more objective benchmarking in order to make a fairer comparison of different
445 tools.

446 **How reliable are the obtained results?**

447 Currently, there is no benchmark evaluation of tools for DNA sequence or scRNA-seq data.
448 Here we attempted to make a preliminary examination of the reliability of results obtained
449 from different pA prediction tools, using a matched bulk RNA-seq and 10x Chromium
450 scRNA-seq data of human peripheral blood mononuclear cells (PBMCs) (File S6). We
451 chose representative tools from each category, including DaPars2 [91], TAPAS [41], and
452 Aptardi [99] for bulk RNA-seq, Sierra [44], scAPATrap [45], and SCAPTURE [107] for

453 scRNA-seq, and DeepPASTA [28] for DNA sequences (**Figure 3A** top). We collected a
454 total of 676,424 non-redundant pAs from GENCODE v39, PolyASite 2.0, and PolyA_DB
455 3, which were compiled from 3' seq and can be used as the true reference (Figure 3A
456 bottom). The number of pAs predicted by different tools, even those under the same
457 category, varies greatly (Figure 3B, left). For example, the number of pAs predicted from
458 RNA-seq by TAPAS and DaPars was nearly 8 times and 4 times that of Aptardi. The
459 number of pAs predicted from scRNA-seq by Sierra and scAPATrap is about twice that of
460 SCAPTURE. Of note, scAPATrap can predict pAs for the whole genome including
461 intergenic regions and all the three tools predict a large number of pAs in introns (Figure
462 3B, right). If only 3' UTR regions are considered, the number of pAs predicted by the three
463 scRNA-seq tools is much closer (Figure 3C, left). As most tools only identify pAs in 3'
464 UTR, here we used 3' UTR pAs for subsequent evaluation. Next, we assessed the
465 authenticity of the predicted pAs by checking whether they are supported by annotated pAs
466 in the true reference. The overlap of pAs predicted from RNA-seq with annotated pAs is
467 much lower than that of scRNA-seq (Figure 3C, left). Particularly, the overlap rate between
468 pAs predicted by SCAPTURE and annotated pAs is as high as 96%, which may be because
469 that the DL model embedded in SCAPTURE was trained with annotated pAs. The position
470 of pAs predicted by TAPAS, SCAPTURE, and scAPATrap is much more precise than that
471 by other tools (Figure 3C, right). Further, we examined the consistency of the results
472 predicted by different tools. Generally, the consistency among different tools is very low
473 (Figure 3D). For RNA-seq data, only 289 pAs were identified by all the three tools, whereas
474 the vast majority of pAs were identified exclusively by a single tool (Figure 3D, top). In
475 contrast, the consistency of pAs predicted from scRNA-seq by different tools is relatively
476 higher (Figure 3D, bottom). In addition, we assessed the reliability of predicted pAs by
477 investigating sequence features. The single nucleotide profile around pAs predicted by
478 TAPAS, SCAPTURE, and scAPATrap resembles the general profile [50] (Figure 3E), which
479 is also consistent with the fact that they determine more precise locations for pAs (Figure

480 3C, right). The percentage of AATAAA around pAs predicted by scRNA-seq is much
481 higher than that of bulk RNA-seq (Figure 3F), indicating that predicted pAs from scRNA-
482 seq tend to be more reliable and more accurate than that from bulk RNA-seq. Next, we
483 used the pA prediction tool for DNA sequence, DeepPASTA, to examine how many pAs
484 identified from RNA-seq data are predicted as true solely based on the sequence
485 characteristics. We extracted the upstream and downstream sequences of pAs predicted by
486 RNA-seq tools as the input for DeepPASTA. The proportion of pAs obtained by different
487 tools to be predicted as true by DeepPASTA is not high and varies greatly, ranging from
488 28% to 79% (Figure 3G, top), indicating again the low overlap of pAs predicted by different
489 tools. Considering only positive pAs by DeepPASTA, the percentage of AATAAA and 1-
490 nt variants of different tools increased slightly (Figure 3G, bottom vs., Figure 3F),
491 reflecting that positive pAs confirmed by DeepPASTA are relatively more reliable than
492 negative ones. Finally, we examined predicted pAs of the immunoglobulin M heavy chain
493 (*IGMM*) gene, which was reported to express a secreted form using the proximal pA and
494 the membrane-bound form using the distal one [123]. The proximal pA of *IGHM* has been
495 recently found preferentially used in B cells and plasma cells of COVID-19 samples [107].
496 SCAPTURE and scAPATrap predicted the precise location of both proximal and distal pAs
497 from scRNA-seq data, while Sierra only predicted the proximal one (Figure 3H). TAPAS
498 predicted three pAs from bulk RNA-seq data, of which two perfectly matched the reference
499 pAs in PolyASite 2.0. In contrast, Aptardi failed to predict any pA for this gene and
500 DaPars2 predicted two pAs yet not verified by reference pAs.

501 Although this preliminary benchmark is far from objective or exhaustive to reflect the
502 advantages and disadvantages of different tools, it reveals several potential issues when
503 using the results obtained by different pA prediction methods. First, although a
504 considerable number of pAs were identified by most tools, the overall prediction accuracy
505 and sensitivity of these tools is low (Figure 3C). Our previous comparative study [50] on
506 tools for bulk RNA-seq have also revealed that a considerable number of predicted pAs

507 were not annotated in 3' seq, and the overall prediction accuracy of these tools, even the
508 best one, TAPAS, is not high (40%–60% for human/mouse data). It is still challenging to
509 determine whether a pA not present in prior annotations is false or novel. We anticipate that
510 at least part of predicted pAs that are not overlapping with annotated ones may potentially
511 be true due to that the current pA annotations are still far from complete. Second, the
512 number of pAs identified by different tools, either for bulk RNA-seq or scRNA-seq, varies
513 greatly, and the consensus of results obtained by different tools is limited (Figure 3D). This
514 is also similar to the observation in our previous benchmark that each tool predicts an
515 independent set of pAs and the overlap of results from different tools is extremely low (<
516 7% for human/mouse data) [50]. Third, as some tools incorporate additional information
517 to predict pAs, e.g., prior pAs used by SCAPTURE and poly(A) reads used by scAPATrap,
518 the resolution of pAs predicted by different tools varies greatly (Figures 3C&E). Fourth,
519 21% to 72% of the predicted pAs by different tools were not recognized as true pAs based
520 on their sequence features (Figure 3G). Fifth, although scRNA-seq data suffers from
521 extremely high level of noise and sparsity, prediction results from scRNA-seq seem to be
522 more reliable and consistent than those from bulk RNA-seq (Figures 3C, D &F). However,
523 this is not unexpected because that it may be less challenging to computationally predict
524 pAs from the 3'-tag based scRNA-seq data than the full-length-based bulk RNA-seq data.
525 Still, further benchmark study with more complete prior annotations, diverse datasets, and
526 performance indicators is needed in order to assess the results obtained from different tools
527 more fairly and objectively.

528 Here we try to give some operable suggestions on how to obtain high-confidence pAs.
529 The most straightforward way may be making a consensus set of pAs that are predicted by
530 multiple tools, however, this may result in a relatively small number of pAs due to the
531 limited overlap by different tools. Another way is to obtain the intersection of predicted
532 sites and real sites, using annotated pAs that are manually curated and available in several
533 databases such as PolyASite 2.0 and PolyA_DB 3. However, it should be noted that these

534 annotated data sources were compiled from limited biological samples and species; they
535 are far from complete to cover all real sites especially tissue-specific ones. Similar to our
536 benchmark analysis on RNA-seq and scRNA-seq PBMCs (Figure 3), users can also use
537 data from another omics from similar biological samples, if available, to predict pAs for
538 mutual verification. In addition, since many sequence motifs, e.g., AAUAAA and its
539 variants, have been reported to have a positional preference relative to the pA, it is naturally
540 to examine sequence patterns surrounding each predicted pA to get pAs with explicit
541 poly(A) signals. This is particularly useful for assessing the authenticity of pAs from
542 animals because AAUAAA and its 1-nt variants appeared in > 90% of animal pAs [8]. In
543 contrast, AAUAAA only accounts for < 10% of pAs in plants, therefore it is not practical
544 to validate plant pAs through sequence features. Moreover, the general single nucleoside
545 compositions surrounding pAs in different species have been clearly reported, we can thus
546 inspect the base composition around predicted pAs. Of note, this way is applicable to
547 evaluation of the overall quality of the pAs, while it cannot be used to assess the reliability
548 of a single pA. The movAPA package [56] can be used for most of the above-mentioned
549 quality assessments.

550 **Practical guidelines for choosing appropriate methods**

551 Based on the summary of different methods (Files S1-S4), we attempt to choose
552 representative tools from each category and propose a set of practical guidelines for users
553 (**Table 1**). As methods in different categories use different kinds of data as the input, the
554 choice of the method first depends on the users' own data. For bulk RNA-seq data, the
555 choice of the method should be mainly driven by the availability of pA annotations. For
556 scRNA-seq data, the choice of the method mainly depends on the protocol of the scRNA-
557 seq (e.g., 3' tag or full-length) and the availability of pA annotations. For methods
558 predicting pAs from DNA sequence, the choice of the method should be primarily driven
559 by the algorithm used, deep learning or traditional machine learning. Particularly, for cross-
560 species pA prediction from DNA sequences, users should pay extra attention to whether

561 they need to retrain the model for individual species, which may require users to have
562 certain programming ability. Additionally, several tools are in the form of web servers,
563 providing a portable platform for predicting pAs from DNA sequences for researchers with
564 limited programming ability. Several other factors also affect the choice of methods, such
565 as the availability of the tool or code, the popularity, the ease of use, the clarity of
566 documentation, and the scale of the data. When predicting pAs on a dataset of interest, it
567 is important to further consider two points. First, it is critical that the obtained pAs and/or
568 the downstream results (e.g., differential APA events) are confirmed by multiple pA
569 prediction methods. This is to ensure that the prediction is not biased due to predefined
570 parameter settings or the specific algorithm used in the method. The merit of using different
571 methods is also demonstrated by the benchmark results in previous studies [50, 51] and in
572 this study (Figure 3), which show substantial complementarity between different methods.
573 Second, even if prior pA annotations are available, it can be also beneficial to try out
574 methods that do not rely on prior annotations. When predicted pAs, even a small portion,
575 are confirmed using such a different method, it provides users with additional evidence.

576 **Conclusions and prospects**

577 **Challenges in improving the performance of pA prediction**

578 The field of pA prediction is progressing rapidly, primarily in the aspects of using DL
579 models and predicting pAs at the single-cell resolution. However, the overall accuracy,
580 sensitivity, and specificity of currently available methods remain moderate (Figure 3). The
581 coming flood of extensive sequencing data, especially multi-omics and single-cell data,
582 will provide new opportunities but also demand new computational methods to exploit this
583 new information. Potential challenges of improving the prediction performance include but
584 are not limited to: paucity of annotated pAs covering diverse tissues and species; mis-
585 assemblies caused by the low complexity 3' UTR sequences; mis-alignment of short reads
586 or incomplete sequence coverage near 3' ends; difficulty in capturing pAs in low-
587 expression genes; poor knowledge on primary, secondary or higher structure information

588 of poly(A) signals, particularly in plants; gaps in our knowledge on understanding APA
589 regulators in different omics layers; limited success in integrating the quantitative features
590 from multiple omics layers; lack of transferrable intelligent methods for cross-species
591 prediction; lack of interpretability in models based on deep neural networks; hurdle in
592 constructing negative datasets due to the prevalence of unconventional pAs in CDS and
593 introns; difficulty in identifying multiple pAs anywhere in the transcript; lack of effective
594 algorithms to deal with the extremely high isoform-level dropout rate and noise inherent in
595 scRNA-seq. Furthermore, higher standards for software quality assurance and
596 documentation would help improve the ease of use of these tools and facilitate their
597 application in the broader community. Finally, new algorithms should be designed to cope
598 with ever-increasing amount of different kinds of data, especially the explosion in single-
599 cell data with multi-omics features.

600 **Notes on benchmarking different methods for predicting pAs**

601 Till now, there are few reports on the exhaustive evaluation of computational tools for
602 predicting pAs. Previously, our group benchmarked 11 representative tools for predicting
603 pAs and/or dynamic APA events from RNA-seq [50]. Lately, Shah et al. [51] evaluated five
604 tools for RNA-seq against 3' seq, Iso-Seq, and a full-length RNA-seq method in identifying
605 pAs and quantifying pA usage. However, there is no study to provide an exhaustive
606 evaluation of existing tools for pA prediction from different kinds of data, particularly those
607 tools for scRNA-seq. Here we attempt to give some notes on benchmarking analysis in this
608 field. First, the real pA dataset is very critical for performance evaluation, however, the
609 reference datasets used in different studies are quite different. Therefore, it is imperative to
610 compile reliable reference datasets with uniform standards. In particular, RNA-seq or
611 scRNA-seq data are sample-specific, so the reference pA dataset from matched samples
612 should also be considered. Moreover, due to the paucity of real pA dataset at the single-cell
613 level, possible deviations need to be considered when using real pA data from bulk data for
614 evaluation. For example, pAs exclusively recognized in single cells may be authentic pAs

615 from rare transcripts or rare cells, even though they may not be present in the bulk pA
616 reference. Second, most tools were evaluated using data only in mammals (mainly human
617 and mouse), therefore the scalability of these tools in different species, especially their
618 applicability to plants, needs to be further evaluated. Third, almost all published prediction
619 tools provide their own benchmark pipelines using different datasets, which potentially
620 favors their prediction efficiency. These benchmark protocols might be credible, but may
621 lack objectivity, simplicity, and effectiveness. We have sorted out the data used for
622 performance evaluation in the respective study of each tool in detail (Files S1-S4), which
623 can facilitate researchers to compile more diverse and standard data for objective
624 benchmark in the future. So far, the most widely used datasets for evaluating pA tools for
625 DNA sequences are the PASS dataset [69, 70] of plant species, the ERPIN dataset [23], and
626 DeepGSR dataset [27] of animal species; datasets for RNA-seq are the MAQC dataset [124]
627 and the HEK293 dataset [125]; datasets for scRNA-seq are the 10x human PBMC data and
628 the *Tabula Muris* atlas [126] (Files S1-S4). Moreover, genomic data could be small sample
629 data and large-scale data, it is also necessary to evaluate the performance of different tools
630 under different sizes of data. Fourth, the output format varies among different tools. For
631 example, most tools for DNA sequences generate binary output or probabilities between 0
632 and 1; some tools for RNA-seq or scRNA-seq output potential regions of pA instead of
633 exact pA position. Therefore, how to unify the output of different tools for objective
634 evaluation needs to be carefully considered. Fifth, compared with the benchmark of tools
635 for DNA sequence data, the benchmark for scRNA-seq tools is much less uniform (Files
636 S1-S4). Almost all studies examined the consensus between the identified pAs and
637 annotated pAs, while there is still no commonly used objective evaluation strategies with
638 diverse indicators. Therefore, it is necessary to use a variety of performance indicators (e.g.,
639 sensitivity, specificity, and precision) that are complementary in nature for comprehensive
640 performance evaluation, particularly for the evaluation of the emerging scRNA-seq tools.
641 At the same time, it is also important to simply present an overall ranking of different tools.

642 The last but not least, many tools have parameters that can be adjusted, however, only the
643 default parameters were normally used for evaluation. Therefore, some strategies (e.g., grid
644 search) should be proposed to evaluate the impact of different parameters of a method.

645 **Predicting pAs in non-3' UTRs**

646 With the advance of 3' seq, more and more unconventional pAs located in non-3' UTR
647 regions like intron and CDS were discovered [3, 49, 127]. These non-3' UTR pAs may
648 generate mRNA isoforms encoding distinct proteins or result in the creation of premature
649 stop codons. Intronic polyadenylation has been found associated with cancer through the
650 inactivation of tumor-suppressor genes [95, 128]. The differential use of intronic pAs is a
651 potential indicator for the differential expression of pre-spliced mRNA transcripts, which
652 contributes to detecting newly transcribed genes and ultimately helps estimate the rate and
653 direction of cell differentiation [129]. Till now, almost all computational tools focused on
654 pA prediction in 3' UTRs. Many tools, particularly those for DNA sequences, usually
655 consider random sequences from introns as negative datasets for model training, which
656 would cause some real intronic pAs to be mistakenly regarded as negative instances.
657 Therefore, even for the pA prediction in 3' UTRs, it is necessary to consider the prevalence
658 of unconventional pAs when constructing the negative dataset. Lately, some tools for bulk
659 or single-cell RNA-seq have found a considerable number of pAs in introns. By applying
660 IPAFinder [95] on pan-cancer data from bulk RNA-seq, 490 recurrent dynamically
661 changed intronic pAs were found. Sierra [44] utilized a splice-aware strategy and identified
662 a considerable number of intronic peaks from scRNA-seq, however, the majority of these
663 peaks may be internal priming artifacts as they are proximal to A-rich regions. SCAPTURE
664 [107] also found > 16,000 candidate intronic pAs from 10x PBMC samples, while < 20%
665 pAs were overlapped with known intronic sites and a large number of false positives were
666 present in lowly expressed genes. Therefore, further careful inspection or filtering is critical
667 to obtain true non-3' UTR pAs or new intelligent algorithms are demanded to effectively
668 call non-3' UTR pAs.

669 **Predicting tissue-specific pAs**

670 APA plays a significant role in tissue-specific regulation of gene expression [2, 12].
671 Profiling APA dynamics or differential APA usages under different physiological or
672 pathological conditions has become a routine analysis in most APA studies. Computational
673 prediction of tissue-specific pAs may be an alternative yet cost-effective solution for
674 analyzing tissue-specificity of APA. The pA prediction problem described in this review is
675 essentially a binary classification problem, which aims to distinguish between nucleotide
676 sequences or genomic regions that contain a pA and those do not. Studies are currently in
677 progress to solve the problem of pA quantification, which aims to predict the strength or
678 dominance of a given pA across tissues. Weng et al. [112] and Hafez et al. [130] predicted
679 whether a given pA is tissue-specific or not, whereas they do not tackle the question of
680 alternative choice of APA sites. One way to study tissue specificity of pAs is to explore the
681 differential usage of APA sites in a gene (e.g., proximal and distal pAs). Several tools, such
682 as Conv-Net [73], have been proposed to predict the strength of APAs sites. Leung et al
683 [73] predicted relative dominance of pAs within 3' UTR in human tissues solely based on
684 nucleotide sequences using a DL model. However, these methods only make predictions
685 based on sequence features, while fail to consider sample-specificity and *in vivo* expression.
686 In contrast, many tools for RNA-seq or scRNA-seq can be used for pairwise comparisons
687 between two samples, while they are not very suitable for profiling APA across multiple
688 tissues. Ever-larger RNA-seq or scRNA-seq data comprising of growing number of
689 samples from diverse tissues are increasingly available, which places new demands on
690 developing new methods to efficiently tackle the question of tissue-specificity of APA.

691 **Cross-species prediction of pAs**

692 Traditional ML methods, such as those based on SVM, can hardly adapt to different species,
693 because they used hand-crafted features learnt for a specific species. Although many DL-
694 based tools have been proposed to improve the performance of pA prediction, most tools
695 still need species-specific real pA collection for model training. Consequently, these tools

696 may suffer from high risks of overfitting and are not applicable to species without any prior
697 pA annotations. Therefore, it is a promising direction to design new transferrable
698 algorithms for cross-species pA prediction or to improve the generalizability of existing
699 tools, which allows a well-trained model from a species with rich annotations to be
700 transferred to data from a different species without retraining or prior knowledge.
701 Annotation-assisted methods, compared to methods without using prior annotations,
702 generally ensure higher data quality and achieve better performance, however their
703 application is limited to data from specific species or biological conditions. Collection of
704 more extensive pA annotations from different sources would definitely contribute to
705 predicting novel sites and increasing the coverage of pAs in diverse cell types, biological
706 conditions, and species. Therefore, an alternative solution for predicting pAs in poorly
707 annotated species could be building an elegant model for well-annotated species and then
708 transferring the model to a different but related species, even without an established pA
709 collection. An initial attempt has been made by some existing methods like Poly(A)-DG
710 [122], which extracts shared features from multiple species and can be generalized to the
711 target species without fine-tuning. However, Poly(A)-DG was only tested between four
712 animals. Till now, tools applicable to plants are still limited. It is widely accepted that the
713 sequence conservation in poly(A) signals in plants is very low, where the most dominant
714 AAUAAA only appears in less than 10% of pAs [57]. Our group recently developed a tool
715 called QuantifyPoly(A) [63] to profile genome-wide polyadenylation choices, which found
716 plant pAs generally exhibit higher micro-heterogeneity than animal ones, and UGUA,
717 UAAA and/or AAUAAA are used in a species-dependent manner. Still, more efforts are
718 needed to explore additional motifs and/or higher-order structures associated with plant
719 polyadenylation and more intelligent algorithms are demanded in order to better predict
720 pAs in multiple species.

721 **Predicting pAs by integrating multi-omics data**

722 Poly(A) sites can be derived from different kinds of data. For example, 3' seq has the unique

723 advantage of acquiring high-quality pAs transcriptome-wide, which contributes to a larger
724 compendium of authentic pAs. Third-generation sequencing technologies, such as PacBio
725 sequencing, are powerful in profiling full-length transcriptome, which could provide a
726 more accurate transcriptome annotation. Widely conducted bulk RNA-seq data can be used
727 for capture and quantify pAs of low-abundance transcripts, and the rapid growing scRNA-
728 seq data support the identification of relatively rare transcripts in single cells. In addition
729 to the genome or transcriptome layers, APA modulation has been found associated with
730 other layers of gene regulation, such as nucleosome positioning, transcription rate, DNA
731 methylation, and RNA-binding proteins [2, 131-133]. By integrating multi-omics data,
732 weak signals from one layer can be amplified or noises be reduced to avoid false negative
733 predictions by referring to the complementary information from additional layers. For
734 instance, potential pAs identified from RNA-seq without well-recognized poly(A) signals
735 could be eliminated if there is no evidence in 3' seq or full-length RNA-seq data. Initial
736 attempts have been made for APA analysis using multi-omics data. scUTRquant [108]
737 incorporates a cleavage site atlas established from a mouse full-length Microwell-seq
738 dataset of 400,000 single cells [109] for filtering high-confidence pAs predicted from 3'
739 tag scRNA-seq. Leung, et al. [73] predicted strength of pAs using nucleotide sequences,
740 considering features from additional layers like nucleosome positioning and RNA/binding
741 protein motifs. IntMAP [97] is a unified ML-based framework, which can fine-tune the
742 contributions of RNA-seq and 3' seq data by tailoring the parameter λ . Currently, DL
743 models have been widely used in predicting pAs from DNA sequences. However, in many
744 cases, DL models fail to make accurate prediction, while patterns of RNA-seq coverage
745 provide clear evidence of polyadenylation, and vice versa [111]. Accordingly, several DL-
746 based tools integrated bulk or single-cell RNA-seq with DNA sequences for pA prediction,
747 such as SCAPTURE [107] and Aptardi [99]. It is promising yet challenging to formulate
748 one unified computational framework, especially leveraging the strength of intelligent
749 algorithms, to integrate the quantitative information from multiple omics layers, e.g.,

750 genomic DNA, transcriptome data, methylation data, and chromatin accessibility data, to
751 identify and quantify pAs genome-wide.

752 **Predicting pAs at the single-cell level**

753 With the rapid development of scRNA-seq technology, different tools continue to emerge
754 for pA identification in single cells. Currently, most methods, like scAPA [102], Sierra [44],
755 and MAPPER [106], construct pseudo-bulk RNA-seq data by pooling reads from cells of
756 the same cell cluster (or cell type) to address the high dropout rate and variability inherent
757 in scRNA-seq. Although many of these tools, like scAPA or Sierra, can still quantify the
758 expression of a pA in each cell by counting reads within a poly(A) region, single cell-based
759 quantification may have high noise level and missing values due to biological and technical
760 variance [106]. As such, APA usage is characterized at the cell-cluster resolution rather than
761 the single-cell resolution, which somewhat contradicts the ultimate goal of single-cell
762 sequencing. Moreover, cell clusters or cell types in these studies were inferred by the
763 conventional gene-cell expression profile, consequently, the APA analysis is limited to
764 predefined cell types and the result may be affected by different cell type annotations.
765 Alternatively, scDaPars [46] quantifies single-cell APA usage based on the model for bulk
766 RNA-seq introduced in DaPars [39], and then recovers missing APA usage by leveraging
767 APA information of the same gene in similar cells. Another limitation of most tools for
768 scRNA-seq is that they are only applicable to 3' tag based scRNA-seq like 10x Chromium
769 or CEL-seq. Till now, only scDaPars can be applied to both 3' tag and full-length scRNA-
770 seq, e.g., Smart-seq2. However, although scDaPars is reported to be able to quantify APA
771 usage in individual cells independent of gene expression, pAs were actually predicted from
772 the bulk RNA-seq tool DaPars that was not specifically designed for scRNA-seq. Moreover,
773 it is challenging to identify and verify low-expression pAs in highly sparse scRNA-seq,
774 particularly those in rare cells. In addition, the tsunami of complex scRNA-seq datasets
775 with various tissue sources, batch effects, and library sizes also have brought huge
776 computational and analytical challenges. Therefore, more efforts are needed to develop

777 new methods to address inherent issues in scRNA-seq for establishing a more
778 comprehensive landscape of pAs at the single-gene and single-cell resolution.

779 **CRediT author statement**

780 **Wenbin Ye:** Investigation, Data curation, Visualization, Writing - review & editing. **Qiwei**
781 **Lian:** Data curation, Writing - review & editing. **Congting Ye:** Investigation, Writing -
782 review & editing. **Xiaohui Wu:** Conceptualization, Writing - original draft, Writing -
783 review & editing, Supervision, Project administration, Funding acquisition. All authors
784 read and approved the final manuscript.

785 **Competing interests**

786 The authors have declared no competing interests.

787 **Acknowledgments**

788 This work was supported by the National Natural Science Foundation of China (Grant No.
789 61871463 to XW) and the Natural Science Foundation of Fujian Province of China (Grant
790 No. 2020J01047 to CY).

791 **ORCID**

792 0000-0002-7811-2710 (Wenbin Ye)
793 0000-0003-3366-6127 (Qiwei Lian)
794 0000-0003-4803-2098 (Congting Ye)
795 0000-0003-0356-7785 (Xiaohui Wu)

796 **References**

797 [1] Wu X, Bartel DP. Widespread influence of 3' -End structures on mammalian mRNA processing and
798 stability. *Cell* 2017;169:905 – 17.e11.
799 [2] Tian B, Manley JL. Alternative polyadenylation of mRNA precursors. *Nat Rev Mol Cell Biol* 2017;18:18–
800 30.
801 [3] Di Giannattino DC, Nishida K, Manley JL. Mechanisms and consequences of alternative
802 polyadenylation. *Mol Cell* 2011;43:853–66.
803 [4] Tian B, Manley JL. Alternative cleavage and polyadenylation: the long and short of it. *Trends Biochem
804 Sci* 2013;38:312–20.
805 [5] Wu X, Liu M, Downie B, Liang C, Ji G, Li QQ, et al. Genome-wide landscape of polyadenylation in

806 Arabidopsis provides evidence for extensive alternative polyadenylation. *Proc Natl Acad Sci USA*
807 2011;108:12533–8.

808 [6] Lianoglou S, Garg V, Yang JL, Leslie CS, Mayr C. Ubiquitously transcribed genes use alternative
809 polyadenylation to achieve tissue-specific expression. *Genes Dev* 2013;27:2380–96.

810 [7] Ulitsky I, Shkumatava A, Jan CH, Subtelny AO, Koppstein D, Bell GW, et al. Extensive alternative
811 polyadenylation during zebrafish development. *Genome Res* 2012;22:2054–66.

812 [8] Gruber AJ, Schmidt R, Gruber AR, Martin G, Ghosh S, Belmadani M, et al. A comprehensive analysis of
813 3' end sequencing data sets reveals novel polyadenylation signals and the repressive role of heterogeneous
814 ribonucleoprotein C on cleavage and polyadenylation. *Genome Res* 2016;26:1145–59.

815 [9] Derti A, Garrett-Engele P, MacIsaac KD, Stevens RC, Sriram S, Chen R, et al. A quantitative atlas of
816 polyadenylation in five mammals. *Genome Res* 2012;22:1173–83.

817 [10] Ozsolak F, Kapranov P, Foissac S, Kim SW, Fishilevich E, Monaghan AP, et al. Comprehensive
818 polyadenylation site maps in yeast and human reveal pervasive alternative polyadenylation. *Cell*
819 2010;143:1018–29.

820 [11] Hoque M, Ji Z, Zheng DH, Luo WT, Li WC, You B, et al. Analysis of alternative cleavage and
821 polyadenylation by 3' region extraction and deep sequencing. *Nat Methods* 2013;10:133–9.

822 [12] Gruber AJ, Zavolan M. Alternative cleavage and polyadenylation in health and disease. *Nat Rev Genet*
823 2019;20:599–614.

824 [13] Oktaba K, Zhang W, Lotz Thea S, Jun David J, Lemke Sandra B, Ng Samuel P, et al. ELAV links paused
825 Pol II to alternative polyadenylation in the *Drosophila* nervous system. *Mol Cell* 2015;57:341–8.

826 [14] Blazie SM, Babb C, Wilky H, Rawls A, Park JG, Mangone M. Comparative RNA-Seq analysis reveals
827 pervasive tissue-specific alternative polyadenylation in *Caenorhabditis elegans* intestine and muscles. *BMC*
828 *Biol* 2015;13:4.

829 [15] Berkovits BD, Mayr C. Alternative 3' UTRs act as scaffolds to regulate membrane protein localization.
830 *Nature* 2015;522:363–7.

831 [16] Batra R, Manchanda M, Swanson MS. Global insights into alternative polyadenylation regulation. *RNA*
832 *Biol* 2015;12:597–602.

833 [17] Xia Z, Donehower LA, Cooper TA, Neilson JR, Wheeler DA, Wagner EJ, et al. Dynamic analyses of
834 alternative polyadenylation from RNA-seq reveal a 3' - UTR landscape across seven tumour types. *Nature*
835 *Communications* 2014;5.

836 [18] Han T, Kim JK. Driving glioblastoma growth by alternative polyadenylation. *Cell Res* 2014;24:1023–4.

837 [19] Gupta I, Clauder-Munster S, Klaus B, Jarvelin AI, Aiyar RS, Benes V, et al. Alternative polyadenylation
838 diversifies post-transcriptional regulation by selective RNA-protein interactions. *Mol Syst Biol* 2014;10:719.

839 [20] Gruber AR, Martin G, Muller P, Schmidt A, Gruber AJ, Gumienny R, et al. Global 3' UTR shortening
840 has a limited effect on protein abundance in proliferating T cells. *Nature Communications* 2014;5.

841 [21] Salamov AA, Solovyev VV. Recognition of 3'-processing sites of human mRNA precursors. *Comput.*
842 *Appl. Biosci.* 1997;13:23–8.

843 [22] Akhtar MN, Bukhari SA, Fazal Z, Qamar R, Shahmuradov I. POLYAR, a new computer program for
844 prediction of poly(A) sites in human sequences. *BMC Genomics* 2010;11:646.

845 [23] Legendre M, Gautheret D. Sequence determinants in human polyadenylation site selection. *BMC*
846 *Genomics* 2003;4:7.

847 [24] Cheng Y, Miura RM, Tian B. Prediction of mRNA polyadenylation sites by support vector machine.
848 Bioinformatics 2006;22:2320–5.

849 [25] Magana-Mora A, Kalkatawi M, Bajic VB. Omni-PolyA: a method and tool for accurate recognition of
850 Poly(A) signals in human genomic DNA. BMC Genomics 2017;18:620.

851 [26] Xia Z, Li Y, Zhang B, Li Z, Hu Y, Chen W, et al. DeeReCT-PolyA: a robust and generic deep learning
852 method for PAS identification. Bioinformatics 2019;35:2371–9.

853 [27] Kalkatawi M, Magana-Mora A, Jankovic B, Bajic VB. DeepGSR: an optimized deep-learning structure
854 for the recognition of genomic signals and regions. Bioinformatics 2019;35:1125–32.

855 [28] Arefeen A, Xiao X, Jiang T. DeepPASTA: deep neural network based polyadenylation site analysis.
856 Bioinformatics 2019;35:4577–85.

857 [29] Bogard N, Linder J, Rosenberg AB, Seelig G. A deep neural network for predicting and engineering
858 alternative polyadenylation. Cell 2019;178:91–106.e23.

859 [30] Sherstnev A, Duc C, Cole C, Zacharaki V, Hornyik C, Ozsolak F, et al. Direct sequencing of *Arabidopsis*
860 *thaliana* RNA reveals patterns of cleavage and polyadenylation. Nat Struct Mol Biol 2012;19:845–52.

861 [31] Jan CH, Friedman RC, Ruby JG, Bartel DP. Formation, regulation and evolution of *Caenorhabditis*
862 *elegans* 3'UTRs. Nature 2011;469:97–101.

863 [32] Harrison PF, Powell DR, Clancy JL, Preiss T, Boag PR, Traven A, et al. PAT-seq: a method to study the
864 integration of 3' -UTR dynamics with gene expression in the eukaryotic transcriptome. RNA 2015;21:1502–
865 10.

866 [33] Park JE, Yi H, Kim Y, Chang H, Kim VN. Regulation of poly(A) tail and translation during the somatic
867 cell cycle. Mol Cell 2016;62:462–71.

868 [34] Chang H, Lim J, Ha M, Kim VN. TAIL-seq: genome-wide determination of poly(A) tail length and 3'
869 end modifications. Mol Cell 2014;53:1044 – 52.

870 [35] Shi Y. Alternative polyadenylation: New insights from global analyses. RNA 2012;18:2105–17.

871 [36] Elkon R, Ugalde AP, Agami R. Alternative cleavage and polyadenylation: extent, regulation and function.
872 Nat Rev Genet 2013;14:496–506.

873 [37] Ji G, Guan J, Zeng Y, Li QQ, Wu X. Genome-wide identification and predictive modeling of
874 polyadenylation sites in eukaryotes. Briefings Bioinf 2015;16:304–13.

875 [38] Ha KCH, Blencowe BJ, Morris Q. QAPA: a new method for the systematic analysis of alternative
876 polyadenylation from RNA-seq data. Genome Biol 2018;19:45.

877 [39] Xia Z, Donehower LA, Cooper TA, Neilson JR, Wheeler DA, Wagner EJ, et al. Dynamic analyses of
878 alternative polyadenylation from RNA-seq reveal a 3'-UTR landscape across seven tumour types. Nat
879 Commun 2014;5:5274–86.

880 [40] Ye C, Long Y, Ji G, Li QQ, Wu X. APAtrap: identification and quantification of alternative
881 polyadenylation sites from RNA-seq data. Bioinformatics 2018;34:1841–9.

882 [41] Arefeen A, Liu J, Xiao X, Jiang T. TAPAS: tool for alternative polyadenylation site analysis.
883 Bioinformatics 2018;34:2521–9.

884 [42] Hashimshony T, Wagner F, Sher N, Yanai I. CEL-Seq: single-cell RNA-seq by multiplexed linear
885 amplification. Cell Rep 2012;2:666–73.

886 [43] Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital
887 transcriptional profiling of single cells. Nat Commun 2017;8:14049.

888 [44] Patrick R, Humphreys DT, Janbandhu V, Oshlack A, Ho JWK, Harvey RP, et al. Sierra: discovery of
889 differential transcript usage from polyA-captured single-cell RNA-seq data. *Genome Biol* 2020;21:167.

890 [45] Wu X, Liu T, Ye C, Ye W, Ji G. scAPATrap: identification and quantification of alternative
891 polyadenylation sites from single-cell RNA-seq data. *Brief Bioinform* 2021;22.

892 [46] Gao Y, Li L, Amos CI, Li W. Analysis of alternative polyadenylation from single-cell RNA-seq using
893 scDaPars reveals cell subpopulations invisible to gene expression. *Genome Res* 2021;31:1856–66.

894 [47] Szkop KJ, Nobeli I. Untranslated parts of genes interpreted: making heads or tails of high-throughput
895 transcriptomic data via computational methods: computational methods to discover and quantify isoforms
896 with alternative untranslated regions. *Bioessays* 2017;39.

897 [48] Yeh HS, Zhang W, Yong J. Analyses of alternative polyadenylation: from old school biochemistry to
898 high-throughput technologies. *BMB Rep* 2017;50:201–7.

899 [49] Chen W, Jia Q, Song Y, Fu H, Wei G, Ni T. Alternative polyadenylation: methods, findings, and impacts.
900 *Genom Proteom Bioinf* 2017;15:287–300.

901 [50] Chen M, Ji G, Fu H, Lin Q, Ye C, Ye W, et al. A survey on identification and quantification of alternative
902 polyadenylation sites from RNA-seq data. *Briefings Bioinf* 2020;21:1261–76.

903 [51] Shah A, Mittleman BE, Gilad Y, Li YI. Benchmarking sequencing methods and tools that facilitate the
904 study of alternative polyadenylation. *Genome Biol* 2021;22:291.

905 [52] Ye C, Lin J, Li QQ. Discovery of alternative polyadenylation dynamics from single cell types. *Comput
906 Struct Biotechnol J* 2020;18:1012–9.

907 [53] Zhang Y, Liu L, Qiu Q, Zhou Q, Ding J, Lu Y, et al. Alternative polyadenylation: methods, mechanism,
908 function, and role in cancer. *J Exp Clin Cancer Res* 2021;40:51.

909 [54] Kandhari N, Kraupner-Taylor CA, Harrison PF, Powell DR, Beilharz TH. The detection and
910 bioinformatic analysis of alternative 3' UTR isoforms as potential cancer biomarkers. *Int J Mol Sci*
911 2021;22:5322.

912 [55] Wang R, Tian B. APAlyzer: a bioinformatic package for analysis of alternative polyadenylation isoforms.
913 *Bioinformatics* 2020.

914 [56] Ye W, Liu T, Fu H, Ye C, Ji G, Wu X. movAPA: Modeling and visualization of dynamics of alternative
915 polyadenylation across biological samples. *Bioinformatics* 2021;37:2470–2.

916 [57] Tian B, Gruber JH. Signals for pre-mRNA cleavage and polyadenylation. *Wiley Interdiscip Rev: RNA*
917 2012;3:385–96.

918 [58] Tabaska JE, Zhang MQ. Detection of polyadenylation signals in human DNA sequences. *Gene*
919 1999;231:77–86.

920 [59] Ahmed F, Kumar M, Raghava GPS. Prediction of polyadenylation signals in human DNA sequences
921 using nucleotide frequencies. *In Silico Biol* 2009;9:135–48.

922 [60] Ji G, Wu X, Shen Y, Huang J, Li QQ. A classification-based prediction model of messenger RNA
923 polyadenylation sites. *J Theor Biol* 2010;265:287–96.

924 [61] Kalkatawi M, Rangkuti F, Schramm M, Jankovic BR, Kamau A, Chowdhary R, et al. Dragon PolyA
925 Spotter: predictor of poly(A) motifs within human genomic DNA sequences. *Bioinformatics* 2012;28:127–
926 9.

927 [62] Wu X, Ji G, Zeng Y. In silico prediction of mRNA poly(A) sites in *Chlamydomonas reinhardtii*. *Mol
928 Genet Genomics* 2012;287:895–907.

929 [63] Ye C, Zhao D, Ye W, Wu X, Ji G, Li QQ, et al. QuantifyPoly(A): reshaping alternative polyadenylation
930 landscapes of eukaryotes with weighted density peak clustering. *Briefings Bioinf* 2021.

931 [64] Liu H, Han H, Li J, Wong L. An in-silico method for prediction of polyadenylation signals in human
932 sequences. *Genome Informatics* 2003;14:84–93.

933 [65] Chang TH, Wu LC, Chen YT, Huang HD, Liu BJ, Cheng KF, et al. Characterization and prediction of
934 mRNA polyadenylation sites in human genes. *Med Biol Eng Comput* 2011;1–10.

935 [66] Xie B, Jankovic BR, Bajic VB, Song L, Gao X. Poly(A) motif prediction using spectral latent features
936 from human DNA sequences. *Bioinformatics* 2013;29:i316–i25.

937 [67] Graber JH, McAllister GD, Smith TF. Probabilistic prediction of *Saccharomyces cerevisiae* mRNA 3'
938 -processing sites. *Nucleic Acids Res* 2002;30:1851 – 8.

939 [68] Hajarnavis A, Korf I, Durbin R. A probabilistic model of 3' end formation in *Caenorhabditis elegans*.
940 *Nucleic Acids Res* 2004;32:3392–9.

941 [69] Ji G, Zheng J, Shen Y, Wu X, Jiang R, Lin Y, et al. Predictive modeling of plant messenger RNA
942 polyadenylation sites. *BMC Bioinf* 2007;8.

943 [70] Shen Y, Ji G, Haas BJ, Wu X, Zheng J, Reese GJ, et al. Genome level analysis of rice mRNA 3'-end
944 processing signals and alternative polyadenylation. *Nucleic Acids Res* 2008;36:3150–61.

945 [71] Ji G, Li L, Li QQ, Wu X, Fu J, Chen G, et al. PASPA: a web server for mRNA poly(A) site predictions
946 in plants and algae. *Bioinformatics* 2015;31:1671–3.

947 [72] Gao X, Zhang J, Wei Z, Hakonarson H. DeepPolyA: a convolutional neural network approach for
948 polyadenylation site prediction. *IEEE Access* 2018;6:24340–9.

949 [73] Leung MKK, Delong A, Frey BJ. Inference of the human polyadenylation code. *Bioinformatics*
950 2018;34:2889–98.

951 [74] Albalawi F, Chahid A, Guo X, Albaradei S, Magana-Mora A, Jankovic BR, et al. Hybrid model for
952 efficient prediction of poly(A) signals in human genomic DNA. *Methods* 2019;166:31–9.

953 [75] Yu H, Dai Z. SANPolyA: a deep learning method for identifying poly(A) signals. *Bioinformatics*
954 2020;36:2393–400.

955 [76] Guo Y, Zhou D, Li W, Cao J, Nie R, Xiong L, et al. Identifying polyadenylation signals with biological
956 embedding via self-attentive gated convolutional highway networks. *Appl Soft Comput* 2021;103:107133.

957 [77] Pickrell JK, Marioni JC, Pai AA, Degner JF, Engelhardt BE, Nkadori E, et al. Understanding
958 mechanisms underlying human gene expression variation with RNA sequencing. *Nature* 2010;464:768–72.

959 [78] Birol I, Raymond A, Chiu R, Nip KM, Jackman SD, Kreitzman M, et al. Kleat: cleavage site analysis of
960 transcriptomes. *Pac Symp Biocomput* 2015:347–58.

961 [79] Bonfert T, Friedel CC. Prediction of poly(A) sites by poly(A) read mapping. *PLoS One*
962 2017;12:e0170914.

963 [80] Campbell MA, Haas BJ, Hamilton JP, Mount SM, Buell CR. Comprehensive analysis of alternative
964 splicing in rice and comparative analyses with *Arabidopsis*. *BMC Genomics* 2006;7:327.

965 [81] Guttman M, Garber M, Levin JZ, Donaghey J, Robinson J, Adiconis X, et al. Ab initio reconstruction of
966 cell type-specific transcriptomes in mouse reveals the conserved multi-exonic structure of lincRNAs. *Nat
967 Biotechnol* 2010;28:503–10.

968 [82] Le Pera L, Mazzapoda M, Tramontano A. 3USS: a web server for detecting alternative 3'UTRs from
969 RNA-seq experiments. *Bioinformatics* 2015;31:1845–7.

970 [83] Huang Z, Teeling EC. ExUTR: a novel pipeline for large-scale prediction of 3' -UTR sequences from
971 NGS data. *BMC Genomics* 2017;18:847.

972 [84] Wilkening S, Pelechano V, Jarvelin AI, Tekkedil MM, Anders S, Benes V, et al. An efficient method for
973 genome-wide polyadenylation site mapping and RNA quantification. *Nucleic Acids Res* 2013;41:e65-e.

974 [85] Wang R, Nambiar R, Zheng D, Tian B. PolyA_DB 3 catalogs cleavage and polyadenylation sites
975 identified by deep sequencing in multiple genomes. *Nucleic Acids Res* 2018;46:D315-d9.

976 [86] Zhu S, Ye W, Ye L, Fu H, Ye C, Xiao X, et al. PlantAPAdb: a comprehensive database for alternative
977 polyadenylation sites in plants. *Plant Physiol* 2020;182:228–42.

978 [87] Gruber AJ, Schmidt R, Ghosh S, Martin G, Gruber AR, van Nimwegen E, et al. Discovery of
979 physiological and cancer-related regulators of 3' UTR processing with KAPAC. *Genome Biol* 2018;19:44.

980 [88] Fahmi NA, Chang J-W, Nassereddeen H, Ahmed KT, Fan D, Yong J, et al. APA-Scan: Detection and
981 Visualization of 3'-UTR APA with RNA-seq and 3'-end-seq Data. *bioRxiv* 2020:2020.02.16.951657.

982 [89] Kim M, You BH, Nam JW. Global estimation of the 3' untranslated region landscape using RNA
983 sequencing. *Methods* 2015;83:111–7.

984 [90] Shenker S, Miura P, Sanfilippo P, Lai EC. IsoSCM: improved and alternative 3' UTR annotation using
985 multiple change-point inference. *RNA* 2015;21:14 – 27.

986 [91] Li L, Huang K-L, Gao Y, Cui Y, Wang G, Elrod ND, et al. An atlas of alternative polyadenylation
987 quantitative trait loci contributing to complex trait and disease heritability. *Nat Genet* 2021;53:994–1005.

988 [92] Feng X, Li L, Wagner EJ, Li W. TC3A: the cancer 3' UTR atlas. *Nucleic Acids Res* 2018;46:D1027-d30.

989 [93] Zhang J, Wei Z. An empirical Bayes change-point model for identifying 3' and 5' alternative splicing by
990 next-generation RNA sequencing. *Bioinformatics* 2016;32:1823–31.

991 [94] Cass AA, Xiao X. mountainClimber Identifies alternative transcription start and polyadenylation sites
992 in RNA-Seq. *Cell Syst* 2019;9:393–400.e6.

993 [95] Zhao Z, Xu Q, Wei R, Wang W, Ding D, Yang Y, et al. Cancer-associated dynamics and potential
994 regulators of intronic polyadenylation revealed by IPAfinder using standard RNA-seq data. *Genome Res*
995 2021;31:2095–106.

996 [96] Gruber AJ, Gypas F, Riba A, Schmidt R, Zavolan M. Terminal exon characterization with TECtool
997 reveals an abundance of cell-specific isoforms. *Nat Methods* 2018;15:832–6.

998 [97] Chang JW, Zhang W, Yeh HS, Park M, Yao C, Shi Y, et al. An integrative model for alternative
999 polyadenylation, IntMAP, delineates mTOR-modulated endoplasmic reticulum stress response. *Nucleic
1000 Acids Res* 2018;46:5996–6008.

1001 [98] Yang C, Li C, Nip KM, Warren RL, Birol I. Terminitor: cleavage site prediction using deep learning
1002 models. *bioRxiv* 2020:710699.

1003 [99] Lusk R, Stene E, Banaei-Kashani F, Tabakoff B, Kechris K, Saba LM. Aptardi predicts polyadenylation
1004 sites in sample-specific transcriptomes using high-throughput RNA sequencing and DNA sequence. *Nat
1005 Commun* 2021;12:1652.

1006 [100] Ziegenhain C, Vieth B, Parekh S, Reinius B, Guillaumet-Adkins A, Smets M, et al. Comparative
1007 analysis of single-cell RNA sequencing methods. *Mol Cell* 2017;65:631–43.e4.

1008 [101] Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al. Highly parallel genome-wide
1009 expression profiling of individual cells using Nanoliter Droplets. *Cell* 2015;161:1202–14.

1010 [102] Shulman ED, Elkon R. Cell-type-specific analysis of alternative polyadenylation using single-cell

1011 transcriptomics data. *Nucleic Acids Res* 2019;47:10027–39.

1012 [103] Yang Y, Paul A, Bach TN, Huang ZJ, Zhang MQ. Single-cell alternative polyadenylation analysis
1013 delineates GABAergic neuron types. *BMC Biol* 2021;19:144.

1014 [104] Zhou R, Xiao X, He P, Zhao Y, Xu M, Zheng X, et al. SCAPE: a mixture model revealing single-cell
1015 polyadenylation diversity and cellular dynamics during cell differentiation and reprogramming. *Nucleic
1016 Acids Res* 2022.

1017 [105] Meyer E, Dehghannasiri R, Chaung K, Salzman J. ReadZS detects developmentally regulated RNA
1018 processing programs in single cell RNA-seq and defines subpopulations independent of gene expression.
1019 *bioRxiv* 2021:462469.

1020 [106] Li WV, Zheng D, Wang R, Tian B. MAAPER: model-based analysis of alternative polyadenylation
1021 using 3' end-linked reads. *Genome Biol* 2021;22:222.

1022 [107] Li GW, Nan F, Yuan GH, Liu CX, Liu X, Chen LL, et al. SCAPTURE: a deep learning-embedded
1023 pipeline that captures polyadenylation information from 3' tag-based RNA-seq of single cells. *Genome Biol*
1024 2021;22:221.

1025 [108] Fansler MM, Zhen G, Mayr C. Quantification of alternative 3' UTR isoforms from single cell RNA-
1026 seq data with scUTRquant. *bioRxiv* 2021:2021.11.22.469635.

1027 [109] Han X, Wang R, Zhou Y, Fei L, Sun H, Lai S, et al. Mapping the mouse cell atlas by Microwell-Seq.
1028 *Cell* 2018;172:1091–107.e17.

1029 [110] Levin M, Zalts H, Mostov N, Hashimshony T, Yanai I. Gene expression dynamics are a proxy for
1030 selective pressures on alternatively polyadenylated isoforms. *Nucleic Acids Res* 2020;48:5926–38.

1031 [111] Li Z, Li Y, Zhang B, Li Y, Long Y, Zhou J, et al. DeeReCT-APA: prediction of alternative
1032 polyadenylation site usage through deep learning. *Genom Proteom Bioinf* 2021.

1033 [112] Weng L, Li Y, Xie X, Shi Y. Poly(A) code analyses reveal key determinants for tissue-specific mRNA
1034 alternative polyadenylation. *RNA* 2016;19:19.

1035 [113] Ji G, Chen M, Ye W, Zhu S, Ye C, Su Y, et al. TSAPA: identification of tissue-specific alternative
1036 polyadenylation sites in plants. *Bioinformatics* 2018;34:2123–5.

1037 [114] Ye C, Zhou Q, Hong Y, Li QQ. Role of alternative polyadenylation dynamics in acute myeloid
1038 leukaemia at single-cell resolution. *RNA Biol* 2019;16:785–97.

1039 [115] Lu J, Bushel PR. Dynamic expression of 3' UTRs revealed by Poisson hidden Markov modeling of
1040 RNA-Seq: implications in gene expression profiling. *Gene* 2013;527:616–23.

1041 [116] Wang W, Wei Z, Li H. A change-point model for identifying 3'UTR switching by next-generation RNA
1042 sequencing. *Bioinformatics* 2014;30:2162–70.

1043 [117] Katz Y, Wang ET, Airoldi EM, Burge CB. Analysis and design of RNA sequencing experiments for
1044 identifying isoform regulation. *Nat Methods* 2010;7:1009–15.

1045 [118] Grassi E, Mariella E, Lembo A, Molineris I, Provero P. Roar: detecting alternative polyadenylation
1046 with standard mRNA sequencing libraries. *BMC Bioinf* 2016;17:423.

1047 [119] Burri D, Zavolan M. Shortening of 3' UTRs in most cell types composing tumor tissues implicates
1048 alternative polyadenylation in protein metabolism. *RNA* 2021.

1049 [120] Bai Y, Qin Y, Fan Z, Morrison RM, Nam K, Zarour HM, et al. scMAPA: Identification of cell-type-
1050 specific alternative polyadenylation in complex tissues. *GigaScience* 2022;11.

1051 [121] Ye C, Zhou Q, Wu X, Yu C, Ji G, Saban DR, et al. scDAPA: detection and visualization of dynamic

1052 alternative polyadenylation from single cell RNA-seq data. *Bioinformatics* 2019;36:1262–4.

1053 [122] Zheng Y, Wang H, Zhang Y, Gao X, Xing EP, Xu M. Poly(A)-DG: A deep-learning-based domain

1054 generalization method to identify cross-species poly(A) signal without prior knowledge from target species.

1055 *PLoS Comput Biol* 2020;16:e1008297.

1056 [123] Singh I, Lee S-H, Sperling AS, Samur MK, Tai Y-T, Fulciniti M, et al. Widespread intronic

1057 polyadenylation diversifies immune cell transcriptomes. *Nat Commun* 2018;9:1716.

1058 [124] Bullard JH, Purdom E, Hansen KD, Dudoit S. Evaluation of statistical methods for normalization and

1059 differential expression in mRNA-Seq experiments. *BMC Bioinf* 2010;11:94.

1060 [125] Liu N, Dai Q, Zheng G, He C, Parisien M, Pan T. N(6)-methyladenosine-dependent RNA structural

1061 switches regulate RNA-protein interactions. *Nature* 2015;518:560–4.

1062 [126] Schaum N, Karkanias J, Neff NF, May AP, Quake SR, Wyss-Coray T, et al. Single-cell transcriptomics

1063 of 20 mouse organs creates a *Tabula Muris*. *Nature* 2018;562:367–72.

1064 [127] de Lorenzo L, Sorenson R, Bailey-Serres J, Hunt AG. Noncanonical Alternative Polyadenylation

1065 Contributes to Gene Regulation in Response to Hypoxia. *The Plant Cell* 2017;29:1262–77.

1066 [128] Lee SH, Singh I, Tisdale S, Abdel-Wahab O, Leslie CS, Mayr C. Widespread intronic polyadenylation

1067 inactivates tumour suppressor genes in leukaemia. *Nature* 2018;561:127–31.

1068 [129] La Manno G, Soldatov R, Zeisel A, Braun E, Hochgerner H, Petukhov V, et al. RNA velocity of single

1069 cells. *Nature* 2018;560:494–8.

1070 [130] Hafez D, Ni T, Mukherjee S, Zhu J, Ohler U. Genome-wide identification and predictive modeling of

1071 tissue-specific alternative polyadenylation. *Bioinformatics* 2013;29:i108–i16.

1072 [131] Neve J, Patel R, Wang Z, Louey A, Furger AM. Cleavage and polyadenylation: Ending the message

1073 expands gene regulation. *RNA Biol* 2017;14:1–26.

1074 [132] Mayr C. Regulation by 3'-Untranslated Regions. *Annu Rev Genet* 2017;51:171–94.

1075 [133] MacDonald CC. Tissue-specific mechanisms of alternative polyadenylation: Testis, brain, and beyond

1076 (2018 update). *Wiley Interdiscip Rev RNA* 2019:e1526.

1077

1078

1079 **Figure legends**

1080 **Figure 1 Schematic of computational approaches for predicting pAs from different**
1081 **kinds of sequencing data**

1082 **A.** Predicting pAs from DNA sequences based on traditional ML models. **B.** Predicting
1083 pAs from DNA sequences based on DL models. **C.** Identifying pAs from 3' seq data. **D.**
1084 Predicting pAs from bulk RNA-seq data. **E.** Predicting pAs from single-cell RNA-seq data.
1085 Some representative methods are listed in the text box on the right. pA, poly(A) site; ML,
1086 machine learning; T, true; F, false; DL, deep learning; scRNA-seq, single-cell RNA-seq.

1087

1088 **Figure 2 Landscape of computational approaches for predicting pA from DNA**
1089 **sequences, bulk RNA-seq, and single-cell RNA-seq over time**

1090 scRNA-seq, single-cell RNA-seq; DL, deep learning; LDF, linear discriminant function;
1091 QDF, quadratic discriminant function; HMM, hidden Markov model; SVM, support vector
1092 machine; BN, Bayesian network; ANN, artificial neural network; PM, probabilistic model;
1093 RF, random forest; CC, combined classifier; CP, change point; AN, annotation-based
1094 method; TA, transcript assembly; PK, peak calling; TL, transcript assembly and read
1095 linking.

1096

1097 **Figure 3 Comparison of representative tools for predicting pAs from a matched**
1098 **bulk RNA-seq and scRNA-seq data of human PBMCs**

1099 **A.** Schematic of the benchmark (top) and the collection of reference pAs from GENCODE
1100 v39, PolyASite 2.0, and PolyA_DB 3 (bottom). **B.** Number of pAs obtained by different
1101 tools (left) and distributions of pAs in different genomic regions (right). **C.** Overlap of 3'
1102 UTR pAs predicted by different tools with reference pAs (left) and distributions of distance
1103 from predicted 3' UTR pAs to reference pAs (right). **D.** Overlap of 3' UTR pAs predicted
1104 by different tools from RNA-seq data (top) and scRNA-seq data (bottom). **E.** Single
1105 nucleotide profile around 3' UTR pAs predicted by different tools. For each tool, the

1106 sequence logo of the most dominant motif around the pA identified by DREME was also
1107 shown. **F.** Number of occurrences of AATAAA and 1-nt variants around pAs predicted by
1108 different tools. **G.** The proportion of pAs obtained by different tools to be predicted as
1109 positive or negative by DeepPASTA (top) and the number of occurrences of AATAAA and
1110 1-nt variants around positive pAs (bottom). The upstream and downstream sequences of
1111 pAs predicted by each RNA-seq tool were extracted as the input for DeepPASTA. **H.**
1112 Predicted 3' UTR pAs by different tools for the *IGHM* gene. Tracks from top to the bottom
1113 are gene model, reference pAs from three databases, read coverage from bulk RNA-seq,
1114 predicted pAs from bulk RNA-seq, read coverage for each cell type of scRNA-seq, and
1115 predicted pAs from scRNA-seq. The red triangles on the chromosome strip highlight the
1116 two representative pAs of *IGHM*.

1117 **Tables**

1118 **Table 1 Recommended tools for predicting pAs from DNA sequences, bulk RNA-**
1119 **seq, and single-cell RNA-seq**

1120 **Supplementary material**

1121 **File S1 List of methods for predicting poly(A) sites or poly(A) signals from DNA**
1122 **sequences based on traditional machine learning models**

1123 **File S2 List of methods for predicting poly(A) sites or poly(A) signals from DNA**
1124 **sequences based on deep learning models**

1125 **File S3 List of methods for predicting poly(A) sites from RNA-seq**

1126 **File S4 List of methods for predicting poly(A) sites from single-cell RNA-seq**

1127 **File S5 List of methods or resources for analysis of alternative polyadenylation**
1128 **rather than prediction of poly(A) sites**

1129 **File S6 Materials and methods used in this study**

1130





