

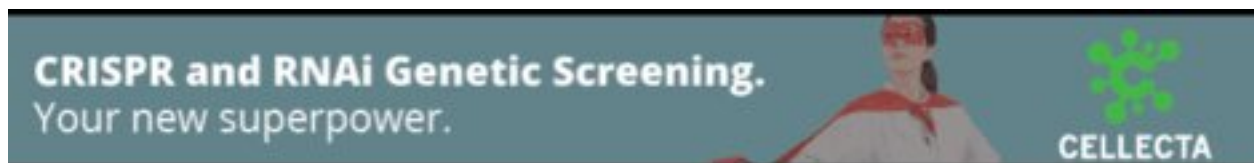


Cancer-associated dynamics and potential regulators of intronic polyadenylation revealed by IPAfinder using standard RNA-seq data

Zhaozhao Zhao, Qiushi Xu, Ran Wei, et al.

Genome Res. published online September 2, 2021
Access the most recent version at doi:[10.1101/gr.271627.120](https://doi.org/10.1101/gr.271627.120)

P<P	Published online September 2, 2021 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



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

Published by Cold Spring Harbor Laboratory Press

1 **Cancer-associated dynamics and potential regulators of intronic**
2 **polyadenylation revealed by IPAFinder using standard RNA-seq data**

3
4 Zhaozhao Zhao¹, Qiushi Xu¹, Ran Wei¹, Weixu Wang¹, Dong Ding¹, Yu Yang¹, Jun
5 Yao¹, Liye Zhang², Yue-Qing Hu³, Gang Wei^{1,*}, Ting Ni^{1,*}.

6
7 ¹State Key Laboratory of Genetic Engineering, Collaborative Innovation Center of
8 Genetics and Development, Human Phenome Institute, School of Life Sciences and
9 Huashan Hospital, Fudan University, Shanghai 200438, P.R. China

10 ²School of Life Science and Technology, ShanghaiTech University, Shanghai 200438,
11 P.R. China

12 ³State Key Laboratory of Genetic Engineering, Institute of Biostatistics, School of Life
13 Sciences, Fudan University, Shanghai 200438, P.R. China

14
15 *Correspondence should be addressed to: G.W. (email: gwei@fudan.edu.cn) and T.N.
16 (email: tingni@fudan.edu.cn)

17
18
19 **Abstract**

20 Intronic polyadenylation (IpA) usually leads to changes in coding region of an mRNA,
21 and its implication in diseases has been recognized, though at its very beginning
22 status. Conveniently and accurately identifying IpA is of great importance for further
23 evaluating its biological significance. Here, we developed IPAFinder, a bioinformatic
24 method for the *de novo* identification of intronic poly(A) sites and their dynamic
25 changes from standard RNA-seq data. Applying IPAFinder to 256 pan-cancer
26 tumor/normal pairs across six tumor types, we discovered 490 recurrent dynamically
27 changed IpA events, some of which are novel and derived from cancer-associated
28 genes such as *TSC1*, *SPERD2*, and *CCND2*. Furthermore, IPAFinder revealed that
29 IpA could be regulated by factors related to splicing and m⁶A modification. In summary,
30 IPAFinder enables the global discovery and characterization of biologically regulated
31 IpA with standard RNA-seq data and should reveal the biological significance of IpA in
32 various processes.

33
34
35

1

2 **Introduction**

3 Alternative polyadenylation (APA) of mRNA is a widespread phenomenon in diverse
4 species, serving as an important contributor to transcriptome diversity (Elkon et al.
5 2013; Tian and Manley 2017). There are different types of APA based on the location
6 of the polyadenylation (pA) site in an mRNA, such as UTR-APA (3' untranslated
7 region APA), CR-APA (coding region APA), and IpA (intronic polyadenylation) (Tian
8 and Manley 2017). UTR-APA changes the length of 3' UTR, thereby altering RNA
9 stability, translation efficiency, RNA localization, or even protein localization (Elkon et
10 al. 2013; Berkovits and Mayr 2015; Tian and Manley 2017). Both CR-APA and IpA
11 introduce a premature termination signal and lead to changes in either the coding
12 sequence or the 3' UTR of the corresponding mRNA (Tian et al. 2007). While
13 UTR-APA is widespread and involved in diverse biological processes (Mayr and Bartel
14 2009; Chen et al. 2018), CR-APA and IpA are less prevalent and their biological
15 functions are not well understood. Recent studies have begun to highlight the
16 biological significance of IpA. For example, IpA diversifies immune cell proteomes via
17 loss of the C-terminal domain (Singh et al. 2018). Cancer cells utilize aberrant intronic
18 pA sites more frequently than normal cells, and the partial loss of function of tumor
19 suppressor genes (TSGs) caused by IpA can contribute crucially to tumorigenesis
20 (Lee et al. 2018). Intronic polyadenylation of *Pdgfra* produces a truncated protein that
21 inhibits PDGF signaling and protects mice from fibrosis (Mueller et al. 2016). *CDK12*
22 suppresses IpA as a mode of regulating DNA repair genes, which is conserved in
23 human tumors that contain loss-of-function *CDK12* mutations (Dubbury et al. 2018).
24 These lines of evidence suggest that genome-wide IpA regulation may play a
25 previously underestimated role in diverse biological processes and pathological
26 conditions.

27 Conveniently and accurately identifying genome-wide IpA is of great importance for
28 further evaluating its biological significance and regulatory mechanism. Although
29 direct 3'-end sequencing of mRNA has provided invaluable insight into the global
30 landscape of APA including IpA (Shepard et al. 2011; Hoque et al. 2013; Ni et al.
31 2013), it has not yet been widely adopted as a routine study strategy, and
32 consequently the availability of such data is currently limited. Conversely, standard
33 RNA-seq has been used in a variety of physiological and pathological conditions, and
34 the amount of related data has increased exponentially in the last decade. Some
35 methods such as DaPars (Xia et al. 2014) and QAPA (Ha et al. 2018) that use
36 RNA-seq data to identify UTR-APA have also been established. However, there is a
37 strong need for a bioinformatic method to *de novo* identify IpA and its dynamic

1 changes using standard RNA-seq data.

2 To meet this demand, we developed a novel bioinformatic algorithm called IPAFinder
3 to identify intronic pA sites and directly infer dynamically changed IpA events by
4 comparative analysis on standard RNA-seq data from different conditions.

5

6 **Results**

7 **IPAFinder identifies dynamically changed IpA events**

8 IPAFinder performs *de novo* identification and quantification of IpA events, without the
9 need for any prior poly(A) site annotation. IpA events are usually classified into two
10 groups: composite terminal exon IpA (or composite IpA) and skipped terminal exon
11 IpA (or skipped IpA) (Tian et al. 2007) (**Supplemental Fig. S1A**). Assuming that there
12 is an intronic poly(A) site (IPA site) used in a given intron, IPAFinder models the
13 normalized RNA-seq read coverage profiles at single-nucleotide resolution and
14 identifies the drop in coverage to infer the potential IPA site, as reflected by the lowest
15 ratio of the sum of mean squared error (MSE) in the upstream and downstream
16 segments split by breakpoint and the MSE computed for the entire intron region
17 ($\text{Ratio}_{\text{MSE}}$) (**Fig. 1A**, see **Methods** for details). Such a strategy could detect composite
18 IpA. To detect skipped IpA (or splicing-coupled IpA), IPAFinder recognizes cryptic 3'
19 splice sites by junction-spanning reads and concatenates the preceding exon to the
20 potential terminal exon (**Supplemental Fig. S1B**). IPAFinder can also exclude
21 alternative splicing events such as alternative 5' splice site and cryptic exon activation
22 using junction-spanning reads to remove potential false-positive events when
23 identifying IpA. Finally, the difference in IpA usage between different conditions can be
24 quantified as changes in the Intronic Poly(A) site Usage Index (ΔIPUI), which can
25 identify dynamically changed IpA events. For example, IPAFinder identified *ELP5* with
26 an increased usage of a composite IPA site in two lung cancer types compared with
27 that in matched normal tissues (**Fig. 1B**). In addition, *ANOS1* showed dynamic
28 changes of splicing-coupled IpA in lung cancers (**Fig. 1C**). The canonical
29 polyadenylation signal (PAS) AAUAAA exists upstream both IPA sites (**Supplemental**
30 **Fig. S2**), which supports the authenticity of these intronic pA sites identified by
31 IPAFinder.

32

33 **Evaluation of IPAFinder using simulated RNA-seq data and experimental 3'-seq** 34 **data**

35 To assess the performance of IPAFinder, we first used simulated RNA-seq data to test
36 whether IPAFinder could accurately infer intronic poly(A) sites. We created simulated
37 data of 5000 genes, of which 500 had a composite IPA site, 500 had a skipped IPA

1 site, 500 had a retained intron, 500 had an alternative 5' splice site, and the others
2 were negative controls. IPAFinder could recover ~80% and ~90% IPA sites at a
3 sequencing coverage of 40x and 60x, respectively (**Supplemental Fig. S3A**).
4 Furthermore, IPAFinder could exclude the interference of alternative splicing events
5 such as alternative 5' splice site and intron retention (**Supplemental Fig. S3B**). Next,
6 we evaluated the ability of IPAFinder to detect dynamically changed IpA events. We
7 simulated 3000 genes with different coverage in two conditions. Among these genes,
8 500 had increasing usage of IPA sites ($\Delta\text{IPUI} > 0.1$), 500 had decreasing usage of IPA
9 sites ($\Delta\text{IPUI} < -0.1$), and 2000 served as negative controls ($|\Delta\text{IPUI}| < 0.1$). If a
10 predicted differentially used IPA site is within 100 bp of a predefined differentially used
11 IPA site then the prediction is considered as a true positive (TP), or otherwise a false
12 positive (FP). For replicated samples, IPAFinder could recover about 80%
13 differentially used IPA sites at a sequencing coverage of 40x with a high precision (**Fig.**
14 **1D**). And the performance of IPAFinder improved with the increase of sequencing
15 depth (**Fig. 1D**). For samples without replicates, we used two methods including
16 bootstrapping-based method and Fisher's exact test-based method to statistically
17 assess the significance of difference for each IpA event between two conditions. As
18 shown, these two methods had comparable sensitivity but Fisher's exact test-based
19 method had better precision (**Fig. 1D**). These results indicate that IPAFinder can infer
20 and quantify IPA sites across a very broad range of RNA-seq coverage levels.

21 Next, we compared IPAFinder-identified IpA events with those found by 3'-seq. We
22 analyzed 3'-seq and standard RNA-seq data of normal and malignant B cells from
23 patients with chronic lymphocytic leukemia (CLL) (Lee et al. 2018). In their original
24 analysis, the authors identified 330 recurrent upregulated IpA events through 3'-seq
25 analysis, followed by validation with standard RNA-seq. Inversely, we first predicted
26 recurrent upregulated IpA events by applying IPAFinder to RNA-seq data, and then
27 used 3'-seq data to validate the results. IPAFinder inferred 306 recurrent upregulated
28 IpA events in malignant B cells compared with those in normal ones, ~84% (256) of
29 which were also found by the original analysis (**Fig. 1E; Supplemental Fig. S4A**).
30 Heatmap of these 256 IpA events also showed an overall increase of IPA site usage in
31 CLL samples, as reflected by lower MSE ratios and higher IPUI values
32 (**Supplemental Fig. S4B**), consistent with the results reported by 3'-seq (Lee et al.
33 2018). These data support the overall agreement between the IPAFinder results
34 based on RNA-seq and those based on 3'-seq. We then undertook a close inspection
35 of those IpA events not overlapping between IPAFinder and the original study (Lee et
36 al. 2018). For IpA events specifically identified by IPAFinder, we found that some
37 genes did have upregulated IpA usage in malignant B cells, as exemplified by

1 *PDXDC1* (**Fig. 1F**). The presence of a noncanonical poly(A) signal AAUACA, a PAS
2 variant ranking 8th among 18 known PASs (Gruber et al. 2016; Ha et al. 2018), was
3 observed upstream the predicted intronic pA site (**Supplemental Fig. S5A,B**). In
4 addition, a clear drop in RNA-seq coverage at the pA site, which has been used for
5 IpA validation (Singh et al. 2018), was observed in *PDXDC1*. The reduced usage of
6 downstream exons of the intronic pA site in *PDXDC1* was also detected in both CLL
7 samples (**Fig. 1F**) and other cancer types such as LUAD and LUSC (**Supplemental**
8 **Fig. S5C–E**). These lines of evidence combined to support the existence and
9 potential regulation of IpA event in *PDXDC1*. For IpA events specifically identified by
10 3'-seq, we found example genes such as *SPTBN1*, which had significantly
11 upregulated IpA usage detected by 3'-seq but did not have significantly higher
12 coverage in the upstream region of the intronic pA site compared with that in the
13 downstream region (**Supplemental Fig. S4C**). As such, IPAFinder could not detect it
14 easily. Based on these results, we conclude that IPAFinder can reliably detect
15 dynamic changes of IpA events between different conditions using standard RNA-seq
16 resources.

17

18 **IPAFinder identifies global landscape of dynamic IpA across tumor types**

19 A previous study demonstrated that IpA could inactivate tumor suppressor genes via
20 the upregulation of truncated mRNAs and proteins and thus contribute to
21 tumorigenesis in CLL (Lee et al. 2018). However, it remains unclear how common
22 IpA-mediated upregulation of truncated mRNAs is in cancers. To examine whether it
23 also occurs in other cancer types, we use TCGA (The Cancer Genome Atlas)
24 database (which archives thousands of RNA-seq data derived from multiple cancer
25 types but lacks the 3' end sequencing data) for IpA analysis. We focused our analysis
26 on six tumor types, namely, lung adenocarcinoma (LUAD), lung squamous cell
27 carcinoma (LUSC), head and neck squamous cell carcinoma (HNSC), prostate
28 adenocarcinoma (PRAD), uterine corpus endometrioid carcinoma (UCEC), and
29 bladder urothelial carcinoma (BLCA), each of which has at least 19 tumor/normal
30 pairs (**Supplemental Table S1**). We identified 130–285 genes with significantly and
31 recurrently (occurrence rate > 10%) changed IpA events for each tumor type and
32 found a total of 490 nonredundant IpA events across the six tumor types (**Fig. 2A**;
33 **Supplemental Fig. S6A,B**; **Supplemental Data S1**). Furthermore, 56% of the 490
34 dynamically changed IpA events occurred in at least two tumor types (**Supplemental**
35 **Fig. S6C**), which indicates the presence of mechanisms potentially acting in concert
36 in IpA regulation across tumor types. Consistent with the phenomenon in CLL, global
37 upregulation of intronic polyadenylation was also prevalent in all six cancer types (**Fig.**

1 **2A**). For IPAFinder-identified intronic pA sites, 45.5% are within 50 nt of the annotated
2 ones compiled from RefSeq, UCSC, Ensembl, and PolyASite 2.0 (Herrmann et al.
3 2020) (**Fig. 2B**). There is an approximately 25-fold enrichment of annotated pA sites
4 in IPAFinder predictions compared with the level in random controls. Enrichment of
5 the canonical poly(A) signal AATAAA in the ± 50 bp flanking sequences of
6 IPAFinder-identified intronic pA sites (**Fig. 2C**) further supported the reliability of our
7 method in discovering IpA in the pan-cancer datasets (Bailey et al. 2009).

8 As mentioned above, intronic polyadenylation tends to disrupt the coding region of
9 an mRNA at different degrees depending on the location at which it occurs. We thus
10 evaluated the impact of IpA on gene expression by computing the fraction of retained
11 coding regions for each IpA isoform relative to the full-length annotated coding regions.
12 An overrepresentation of IpA isoforms that lose the majority of their coding regions
13 was observed (**Fig. 2D**). The remaining IpA isoforms showed a relatively uniform
14 distribution along the coding region (**Fig. 2D**). We found that introns with
15 splicing-coupled IpA events were longer than those with composite IpA or no IpA
16 events (**Fig. 2E**), consistent with the previous finding that a large intron size is a
17 determining factor for IpA events (Tian et al. 2007). A typical example is the *AUH* gene
18 with a splicing-coupled IpA event occurring in an extremely long intron (**Fig. 2F**).
19 These results indicate that IPAFinder can reveal the overall landscape of IpA in six
20 cancer types.

21

22 **Cancer-associated genes are regulated by intronic polyadenylation**

23 To evaluate the relevance of dysregulated IpA events during cancer development, we
24 performed domain analysis and survival analysis on IpA-generated truncated proteins.
25 We found that IpA events occurring at different positions of coding regions all had the
26 possibility of generating truncated proteins with important functional impacts. The first
27 example is *TSC1*, which is required to interact with and stabilize *TSC2* as the
28 *TSC1-TSC2* complex, and the GAP domain on *TSC2* hydrolyzes Rheb-GTP to
29 Rheb-GDP, thereby inhibiting the activation of mTOR kinase (Garami et al. 2003; Inoki
30 et al. 2003; Chong-Kopera et al. 2006). We found that *TSC1* IpA was significantly
31 upregulated in LUAD, LUSC, and HNSC compared with that in normal tissues (**Fig.**
32 **3A,B**). The IpA isoform of *TSC1* was predicted to generate a truncated protein that
33 loses the C-terminal coiled-coil domain (**Fig. 3E**), which is necessary for
34 heterodimerizing with *TSC2* (van Slegtenhorst et al. 1998; Yang et al. 2021). Thus,
35 the IpA-derived truncated protein may fail to form a functional *TSC1/2* complex and
36 lead to the aberrant activation of mTOR kinase. To test this possibility, we transiently
37 co-transfected human HEK293T cells with full-length or IpA-truncated HA-tagged

1 TSC1 and FLAG-tagged TSC2. Then, we performed immunoblot analysis to assess
2 the phosphorylation of S6, an indicator of mTOR kinase activation. Compared with
3 full-length *TSC1*-expressing cells, IpA-truncated *TSC1* failed to inhibit the
4 phosphorylation of S6 (**Fig. 3C**), which indicates that the truncated protein templated
5 by the IpA isoform of *TSC1* has impaired function compared with its full-length version.
6 Furthermore, HNSC patients with higher IpA usage in *TSC1* were found to be
7 associated with worse survival (**Fig. 3D**), which was consistent with previous studies
8 indicating that *Tsc1* inactivation could promote tumor progression in mice (Kladney et
9 al. 2010; Sun et al. 2015). The second example is the *SPRED2* gene. *SPRED2*
10 inhibits the MAP kinase pathway by suppressing the phosphorylation and activation of
11 RAF, in which both EVH-1 and SPR domains are essential (Wakioka et al. 2001;
12 Nobuhisa et al. 2004). *SPRED2* increased the usage of intronic poly(A) sites in
13 multiple cancers (LUAD, LUSC, and HNSC) and thus produced more truncated
14 transcripts with low coding potential (resulting in noncoding RNA as predicted by three
15 different algorithms) (**Fig. 3E,F; Supplemental Figs. S7A,S8**). Experimental
16 validation demonstrated that the truncated transcript of *SPRED2* did not produce any
17 protein in both HEK293T and HeLa cells (**Supplemental Fig. S7B**). Consistent with
18 this, overexpression of the IpA isoform of *SPRED2* failed to inhibit cell proliferation
19 while the full-length version of *SPRED2* did so in the human lung cancer cell line
20 NCI-H520 (**Fig. 3G**). HNSC patients with higher IpA usage in *SPRED2* were also
21 associated with worse survival (**Supplemental Fig. S8C**). The third example is
22 *CCND2*, which encodes cyclin D2. Cyclin D2 has been widely implicated in cell cycle
23 transition and cellular transformation, and its overexpression is highly correlated with
24 poor prognosis in various cancers (Takano et al. 1999; Takano et al. 2000). IPAFinder
25 identified that *CCND2* frequently used IpA to produce a new protein isoform in LUAD,
26 LUSC, and HNSC (**Fig. 3H, Supplemental Fig. S9A–C**). This IpA isoform loses the 3'
27 UTR miRNA repression sites (Mayr and Bartel 2009; Yang et al. 2020) as well as the
28 important residue Thr280 (**Fig. 3E**), which can be phosphorylated by GSK3B and
29 render cyclin D2 susceptible to ubiquitin–proteasome-mediated degradation (Mirzaa
30 et al. 2014). Thus, increased IpA usage in *CCND2* likely causes resistance to protein
31 degradation, which results in *CCND2* accumulating to promote cell-cycle progression.
32 In line with this, the presence of *CCND2* IpA was associated with worse survival in
33 LUSC (**Supplemental Fig. S9D**). These three examples demonstrate that IPAFinder
34 can identify functional IpA regulation in cancer-related genes.

35

1 **Intronic polyadenylation can be influenced by factors related to splicing and**
2 **m⁶A modification**

3 The fidelity of RNA splicing is regulated by an orchestration of splicing enhancers and
4 repressors, and multiple RNA-binding proteins (RBPs) can protect the transcriptome
5 from the aberrant exonization of transposable elements (Zarnack et al. 2013; Attig et
6 al. 2018) and modulate cleavage and polyadenylation at poly(A) sites where they bind
7 (Licatalosi et al. 2008; Hilgers et al. 2012). Applying IPAfinder to published RNA-seq
8 datasets generated from concurrent knockdown of *PTBP1* and *PTBP2* (Gueroussov
9 et al. 2015), two splicing factors preferentially binding to CU repeats (Oberstrass et al.
10 2005), we found that *PTBP1/2* deficiency resulted in many more upregulated IpA
11 events than downregulated ones, consistent with the findings of *PTBP1/2* in
12 repressing cryptic exons in the original study (**Fig. 4A,B**). Sequence analysis
13 confirmed the presence of adjacent CU microsatellites around these activating poly(A)
14 sites (**Fig. 4C**), which suggests the direct binding of *PTBP1/2*. We also tested another
15 RBP heterogeneous nuclear ribonucleoprotein C (*HNRNPC*), which has been
16 reported to modulate the processing of pre-mRNA 3'-ends (Gruber et al. 2016).
17 Applying IPAfinder to RNA-seq datasets of HEK293T cells obtained upon the
18 knockdown of this protein (Liu et al. 2015), we found that the loss of *HNRNPC* also led
19 to the widespread upregulation of IpA events and the majority (72.1%) were skipped
20 IpA events (**Fig. 4D,E**). Sequence analysis of the major IpA isoforms showed that the
21 density of (U)₅ tracts, reported to be the binding site for *HNRNPC* (König et al. 2010),
22 was markedly higher around cryptic 3' splice sites whose usage increased upon
23 *HNRNPC* knockdown compared with those without apparent changes in usage (**Fig.**
24 **4F; Supplemental Fig. S10**). Experimental validation supported the upregulation of
25 IpA events upon knockdown of these RBPs (**Supplemental Fig. S11**). Altogether,
26 these results demonstrate that *PTBP1/2* and *HNRNPC* can protect pre-mRNAs from
27 premature cleavage and polyadenylation by inhibiting the usage of IPA sites.

28 To explore other factors regulating intronic polyadenylation, we next applied
29 IPAfinder to RNA-seq data derived from samples with the knockdown of relevant
30 RBPs. *U2AF1* and *U2AF2* are two auxiliary factors for U2 small nuclear RNA, which
31 bind to the AG dinucleotide and polypyrimidine tract of the 3' splice site, respectively,
32 to facilitate splice site recognition (Zamore et al. 1992; Wu et al. 1999). Multiple
33 studies have reported links between splicing and 3'-end processing (Kyburz et al.
34 2006; Millevoi et al. 2006). Thus, we explored whether these two splicing factors could
35 impact intronic pA site usage by analyzing our custom-made RNA-seq data derived
36 from human foreskin fibroblasts (HFFs) with the knockdown of *U2AF1* or *U2AF2*. We
37 found that depletion of *U2AF1* or *U2AF2* globally increased the usage of intronic pA

1 sites (**Fig. 5A–D**), consistent with a previous analysis of 3'-seq data showing that the
2 knockdown of *U2af2* in mouse C2C12 myoblast cells led to the overall upregulation of
3 IpA events (Li et al. 2015). For skipped IPA sites with increased usage upon the
4 knockdown of *U2AF1* or *U2AF2*, the splicing strength of their cryptic 3' splice sites
5 was significantly weaker than that of downstream 3' splice sites (**Supplemental Fig.**
6 **S12A**). Furthermore, intronic pA sites with increased usage showed considerable
7 overlap between *U2AF1*-KD and *U2AF2*-KD samples (**Supplemental Fig. S12B,C**),
8 suggesting the potential coordination of these two factors in regulating IpA. In line with
9 the similarity in IpA level changes between *U2AF1*-KD and *U2AF2*-KD samples, we
10 also found that both *U2AF1*-KD and *U2AF2*-KD could lead to senescence-associated
11 phenotypes at the cellular level (**Supplemental Fig. S13**) (Yao et al. 2020). Although
12 the causal relationship between IpA and senescence deserves further investigation,
13 we did reveal that *U2AF1* and *U2AF2* have a genome-wide effect on intronic poly(A)
14 site selection.

15 N⁶-methyladenosine (m⁶A) has been identified as the most abundant modification
16 that ubiquitously occurs in eukaryotic mRNAs and affects multiple aspects of mRNA
17 metabolism including alternative splicing (Yang et al. 2018). However, whether factors
18 related to m⁶A modification affect IpA is unclear. Applying IPAfinder to RNA-seq data
19 derived from HeLa cells deficient in *YTHDC1*, the only known m⁶A reader in the
20 nucleus that has been reported to regulate splicing (Xiao et al. 2016), we found that
21 *YTHDC1* deficiency led to increased IPA site usage (**Fig. 5E**). Furthermore, IPA sites
22 with increased usage upon *YTHDC1* knockdown had considerable overlap with those
23 upon *SRSF3* knockdown (**Supplemental Fig. S14**), consistent with previous findings
24 that *YTHDC1* and *SRSF3* intersect with each other to regulate mRNA splicing and 3'
25 UTR length (Xiao et al. 2016; Kasowitz et al. 2018). In addition, analyzing RNA-seq
26 data derived from HEK293T cells deficient in *METTL3*, a methyltransferase implicated
27 in placing m⁶A on RNA (Śledź and Jinek 2016), we found that the loss of *METTL3* also
28 changed the IPA site usage, but the numbers of upregulated or downregulated IpA
29 events were relatively equal (**Fig. 5G**). Example genes with changed IpA usage upon
30 knockdown of either *YTHDC1* or *METTL3* are shown (**Fig. 5F,H**). However, it should
31 be kept in mind that it is currently unclear whether the regulation is directly mediated
32 by m⁶A modification, which warrants further investigation. These results indicate that
33 IpA events have a complex regulatory system and IPAfinder is a valuable tool
34 capable of discovering IpA events from RNA-seq datasets of varied sources, which
35 would also facilitate screening for regulators involved in IpA determination.

36

37 **Comparison of methods for analyzing IpA**

1 IPAFinder was inspired by DaPars, which identifies the breakpoint that can best
2 explain the localized read-density change to perform *de novo* identification and
3 quantification of dynamic UTR-APA events using standard RNA-seq (Xia et al. 2014).
4 However, IpA identification is more complicated than UTR-APA and could be
5 interfered by at least three events, including cryptic exon activation, alternative 5'
6 splice site and intron retention. IPAFinder utilizes BAM format as the input file, which
7 contains splice junction information, and considers most of the interference factors to
8 improve the identification of IpA events (**Supplemental Fig. S15**). We also compared
9 IPAFinder with APALyzer, which analyzes intronic polyadenylation by using RNA-seq
10 data based on known poly(A) sites (such as those annotated in the PolyA_DB
11 database) (Wang et al. 2018; Wang and Tian 2020). Applying APALyzer to RNA-seq
12 datasets obtained upon *HNRNPC*-KD, we observed widespread usage changes of
13 intronic poly(A) sites according to their suggested cutoff (P value < 0.05 using t-test
14 for significance analysis) (**Supplemental Fig. S16A**). Intronic poly(A) sites with
15 increased usage upon *HNRNPC*-KD analyzed by IPAFinder and APALyzer had
16 considerable overlap (**Fig. 6A**). However, IPAFinder could identify dynamic IPA sites
17 that were not annotated by the PolyA_DB database, as exemplified by the IPA sites of
18 *RAD52* and *PTBP2* in *HNRNPC*-KD condition (**Fig. 6B**). IPAFinder could infer
19 upstream splice sites by recognizing junction reads and quantify the usage of
20 corresponding skipped IPA sites accurately, as exemplified by the gene *PPP1R12C*
21 (**Fig. 6B**). Sequence analysis showed that (U)₅ tracts existed in the flanking region of
22 these three IPA sites (**Fig. 6C**), which indicated the direct binding of HNRNPC (König
23 et al. 2010). Although IPAFinder and APALyzer have different criteria for calling
24 differential IpA events, they have a relatively consistent trend in quantifying the usage
25 of IPA sites (**Supplemental Fig. S16 B,C**). Overall, IPAFinder is a specialized tool for
26 *de novo* IpA analysis that is distinct from existing methods such as APALyzer. Different
27 tools have their own strengths and weaknesses (**Supplemental Table S2**), and users
28 may need to apply multiple programs in their research to obtain comprehensive and
29 complementary results.

30

31 **Discussion**

32 In this study, we developed IPAFinder, a method for the *de novo* identification of
33 intronic poly(A) sites and dynamically changed IpA events from standard RNA-seq
34 data. Multiple lines of evidence support the reliability of IPAFinder. Applying IPAFinder
35 to 256 pan-cancer tumor/normal pairs across six tumor types archived in TCGA
36 revealed 490 recurrently changed IpA events, among which there were genes with
37 novel IpA regulation, such as *TSC1*, *SPERD2*, and *CCND2*. Furthermore, genes

1 harboring dynamic IpA events were found being enriched in TSGs but not in the
2 oncogenes (**Supplemental Fig. S17A**). In additional tumor samples (without matched
3 normal tissues), we also found well-known TSGs (*NF1*, *PTEN*, and *CDH1*) with
4 increased IPA site usage (**Supplemental Fig. S17B**). Thus, IPAFinder should help to
5 reveal potential IpA events playing roles in diverse physiological and pathological
6 processes by exploiting the huge amount of standard RNA-seq data.

7 RNA-seq data tends to have coverage biases which is more predominant in
8 untranslated region while we can observe this phenomenon in both real and simulated
9 RNA-seq data (**Supplemental Fig. S18**). Many well-developed methods (such as
10 DaPars, APAtrap and PAQR) (Xia et al. 2014; Gruber et al. 2018; Ye et al. 2018),
11 which *de novo* infer internal poly(A) sites in 3' UTR from standard RNA-seq data, are
12 based on the mean squared error model, regardless of potential coverage bias in 3'
13 UTR. It is worth noting that it is difficult to detect poly(A) sites from RNA-seq data at
14 single nucleotide precision, thus some degree of flexibility (e.g. 100 nt for IPAFinder)
15 is used to match predicted poly(A) sites to the annotated ones (**Supplemental Fig.**
16 **S3**). Although 3'-end sequencing strategies such as 3'-seq coupled with dedicated
17 bioinformatic pipelines can identify IPA sites and detect changes in IPA site usage
18 between different conditions, they are less extensively used in diverse biological
19 processes than standard RNA-seq. In addition, standard RNA-seq has multiple
20 advantages in detecting intronic pA sites compared with 3'-seq: 1) RNA-seq covers
21 the whole gene body and thus junction reads can be used to distinguish skipped IPA
22 sites from composite IPA sites. 2) The use of RNA-seq to infer IPA sites can avoid the
23 interference of internal priming according to continuous upstream read coverage for
24 composite IPA sites and junction-spanning reads for skipped IPA sites.

25 U1 snRNP can protect pre-mRNAs from drastic premature termination by cleavage
26 and polyadenylation at cryptic polyadenylation signals in introns (Kaida et al. 2010;
27 Berg et al. 2012). Applying IPAFinder to publicly available RNA-seq data derived from
28 HeLa cells upon treatment of U1 Antisense Morpholino Oligonucleotide (AMO) (Oh et
29 al. 2017), which has been shown to pair efficiently with U1 snRNA and thereby
30 functionally inhibit U1 snRNP, we found that U1 AMO treatment globally increased the
31 usage of IPA sites (**Supplemental Fig. S19**). These data support the ability of
32 IPAFinder in detecting the usage changes of cryptic IPA sites.

33 We also compared bootstrapping-based method with Fisher's exact test-based
34 method by analyzing RNA-seq dataset obtained by knockdown of *HNRPN*C (merge
35 two replicates as one sample). Bootstrapping-based method identified 197
36 significantly upregulated IpA events while Fisher's exact test-based method identified
37 279 upregulated IpA events, and upregulated IpA events identified by these two

1 methods had considerable overlap (**Supplemental Fig. S20A**). Furthermore, we
2 found that bootstrapping-based method is sensitive to lpA events whose usage
3 difference is relatively large (**Supplemental Fig. S20B**), as exemplified by lpA events
4 of genes *ZCCHC4* and *VPS4B* (**Supplemental Fig. S20C**). By comparing with
5 upregulated lpA events identified by DEXSeq method on replicated samples (Anders
6 et al. 2012), we found that Fisher's exact test-based method could identify more
7 upregulated lpA events supported by DEXSeq method than bootstrapping-based
8 method (174 vs 118). And upregulated lpA events identified by Fisher's exact
9 test-based method have higher fraction supported by DEXSeq method than those
10 identified by bootstrapping-based method (62.4% vs 59.9%). Thus, for samples
11 without replicates between two conditions, users could try both these two statistical
12 methods in their research to obtain comprehensive and complementary results.

13 In conclusion, the IPAFinder method should open up a new avenue for discovering
14 lpA events and changes in their usage in numerous biological processes using
15 standard RNA-seq data. This should help to reveal the functional roles of lpA in
16 diverse conditions.

17

18 **Methods**

19 **IPAFinder algorithm**

20 IPAFinder performs *de novo* identification and quantification of dynamically changed
21 lpA events between two conditions, regardless of any prior poly(A) site annotation.
22 Assuming that there is an intronic pA site used in a given intron, there will be a
23 significant drop in RNA-seq read coverage due to polyadenylation processing. Thus,
24 IPAFinder models the normalized RNA-seq read coverage at single-nucleotide
25 resolution and progressively segments the intron region into two regions with distinct
26 mean coverage. This enables inference of the potential intronic poly(A) site, where the
27 squared deviation decreases most from the mean coverage of the intron when
28 dividing the segment into two regions compared with considering it as a single
29 segment. IPAFinder separately calculates the mean squared error (MSE) of read
30 coverage for upstream (MSE_u) and downstream (MSE_d) segments split by every point
31 in the intron region and compares the sum of MSE_u and MSE_d ($MSE_u + MSE_d$) with the
32 MSE computed for the entire intron region (MSE_e). The ratio of the sum of MSE_u and
33 MSE_d to MSE_e is defined as $Ratio_{MSE}$ (**Fig. 1A**) and, if the lowest value of $Ratio_{MSE} \leq$
34 0.5, a cutoff used to infer internal poly(A) site in 3' UTR by a previous study (Gruber et
35 al. 2018), the corresponding breakpoint is considered as a potential intronic poly(A)
36 site. In addition, the mean coverage in the upstream region of the candidate poly(A)
37 site must be higher than that in the downstream region. Alternative splicing events

1 such as alternative 5' splice site may also have similar segmentation breakpoints, so
 2 we exclude those breakpoints where there are splice sites supported by
 3 junction-spanning reads around them. If the given intron has no composite terminal
 4 exon IpA event, IPAFinder next searches for whether it has a skipped terminal exon
 5 IpA event (**Supplemental Fig. S1B**). We regard a splice site in an internal intron as a
 6 cryptic 3' splice site if it is supported by more than 10 splice junction reads or more
 7 than 10% of upstream 5' splice site junction reads. Then, we concatenate the
 8 preceding exon to this potential skipped terminal exon and identify the best
 9 segmentation breakpoint in the newly formed narrowed intron region, as performed
 10 for the composite terminal exon. Alternative splicing events such as cryptic exon
 11 activation are also excluded by recognizing junction-spanning reads.

12 IPAFinder could also detect multiple IPA sites in a single intron. IPAFinder first infers
 13 the breakpoint with the lowest $\text{Ratio}_{\text{MSE}}$ in the entire intron region. If there is an
 14 alternative intronic poly(A) site in the inferred terminal exon, another drop in RNA-seq
 15 read coverage inside the terminal exon will be observed. Similarly, the alternative
 16 used pA site allows the best segmentation of the terminal exon into upstream and
 17 downstream regions with distinct coverage. Therefore, IPAFinder can infer its location
 18 by calculating the ratio of MSE recursively (**Supplemental Fig. S21A**). The
 19 alternative intronic poly(A) sites of *SPRED2* are identified in such a strategy
 20 (**Supplemental Fig. S21B,C**).

21 Once the intronic poly(A) sites have been identified, library size-normalized
 22 expression levels and relative usage of IPA sites are calculated. We define the Intronic
 23 Poly(A) site Usage Index (IPUI) to quantify the relative IpA usage for sample j as
 24 follows:

$$25 \quad \text{IPUI} = \frac{E_{\text{IPA}}^j}{E_{\text{IPA}}^j + E_{\text{FL}}^j} = \frac{E_{\text{IPA}}^j}{E_{\text{CPE}}^j}, \quad (1)$$

26 where E_{IPA}^j , E_{FL}^j , and E_{CPE}^j are the estimated expression levels of IpA isoform,
 27 full-length isoform, and constitutive preceding exon located upstream of the IPA site
 28 for a given sample (j), respectively. In principle, E_{CPE}^j is equal to the sum of E_{IPA}^j and
 29 E_{FL}^j (**Supplemental Fig. S22**).

30 For samples with replicates, in order to detect differential usage of IpA isoform
 31 between two conditions, we examined the difference in relative usage of terminal
 32 exon inferred by IPAFinder. We applied DEXSeq to model the read counts of all exons
 33 across conditions by negative binomial distribution and tested for the significance of

1 an interaction term between exon and condition (**Supplemental Fig. S23**) (Anders et
2 al. 2012). We defined an lpa isoform to be significantly differentially used if its
3 corresponding terminal exon usage is significantly different between two conditions
4 (FDR-adjusted P value < 0.05) and the difference of IPA site usage is more than 0.1
5 ($|\Delta\text{IPUI}| > 0.1$).

6 For samples without replicates, such as paired tumor-normal samples from TCGA,
7 we used Fisher's exact test to infer differential usage of lpa isoform between
8 conditions, which is a similar approach taken by previous methods for detecting 3'
9 UTR APA events (Xia et al. 2014; Guvenek and Tian 2018). We defined an lpa
10 isoform to be significantly differentially used if its FDR-adjusted P value < 0.05 and
11 $|\Delta\text{IPUI}| > 0.1$.

12 We also provided a bootstrapping-based method to statistically assess the
13 significance of difference for each lpa event between two samples without replicates,
14 which is inspired by the SAAP method (Significance Analysis of Alternative
15 Polyadenylation) (Li et al. 2015). Briefly, for an lpa event from two comparing samples,
16 IPUI was first calculated and was called observed IPUI. Then we sampled reads
17 based on the assumption that the relative abundance of each lpa isoform was the
18 same in two samples. Sampling was performed m times (default $m = 20$) to obtain
19 mean and standard deviation of IPUI, which were then used to convert observed IPUI
20 to Z score. False discovery rate (FDR) and q-value were calculated by comparing
21 observed Z (Z_o) of IPUI and expected Z (Z_e) of IPUI for a given Z cutoff value (Z_c).
22 Q-value for an lpa event x is the FDR using the absolute value of its Z_o (Z_{ox}) as Z_c .
23 We used q-value < 0.05 and $|\Delta\text{IPUI}| > 0.1$ to select significantly differential lpa events.
24 We updated our pipeline in GitHub to indicate which option the users should choose
25 for samples with or without replicates.

26

27 **Data download**

28 All the TCGA RNA-seq BAM files for tumor and matched normal samples were
29 obtained from the Genomic Data Commons (GDC) Data Portal
30 (<https://portal.gdc.cancer.gov/>). Here, we processed LUAD, LUSC, HNSC, UCEC,
31 BLCA, and PRAD cancers. Other publicly available raw sequencing data of 3'-seq
32 and RNA-seq, including those derived from normal immune cells (Singh et al. 2018)
33 and malignant B cells (Lee et al. 2018) from patients with chronic lymphocytic
34 leukemia (CLL) (GSE111310 and GSE111793), *PTBP1/2* knockdown in HEK293 cells
35 (Gueroussov et al. 2015) (GSE69656), *HNRNPC* knockdown and *METTL3*
36 knockdown in HEK293T cells (Liu et al. 2015) (GSE56010), *YTHDC1* knockdown and
37 *SRSF3* knockdown in HeLa cells (Xiao et al. 2016) (GSE71095), *U2AF1* knockdown

1 in HFF cells (Yao et al. 2020) (PRJNA565612), and U1 inhibition in HeLa cells (So et
2 al. 2019) (GSE135140), were obtained from NCBI's Gene Expression Omnibus
3 (GEO).

4

5 **RNA-seq and 3'-seq data analyses**

6 Among the raw paired-end reads obtained from RNA-seq experiments, low-quality
7 reads were filtered out, followed by alignment to the human reference genome
8 sequence (UCSC hg38 assembly) using STAR (Dobin et al. 2013) with the default
9 settings. Analysis of the 3'-seq data was performed as described previously (Singh et
10 al. 2018). The peaks were assigned to RefSeq-annotated genes (downloaded on Jan
11 1, 2020). Isoforms with an expressed level of at least 3 TPM (transcripts per million
12 mapped reads) and usage of at least 0.1 in at least one sample were used for
13 subsequent analyses. We only analyzed lpa isoforms of protein-coding genes.

14

15 **Benchmarking of IPAfinder using simulated RNA-seq data**

16 We first generated a synthetic RNA-seq dataset to assess the performance of
17 IPAfinder to infer intronic poly(A) sites from standard RNA-seq data. To simulate the
18 different coverage levels, baseline coverage for each gene was uniformly sampled
19 between 20x and 80x. An "nx" coverage means that an exonic genomic locus is
20 covered by n reads on average. The usage of lpa isoform or alternative splicing
21 isoform was uniformly sampled from a usage range (40%–60%). We also evaluated
22 the ability of IPAfinder to detect dynamically changed lpa events at different levels of
23 sequencing coverage between two conditions. IPUI values for each gene were
24 randomly sampled until the conditions outlined were met. For samples with replicates,
25 three replicates per condition were generated using negative binomial distribution.
26 The R package polyester was applied to simulate paired-end 100-nt reads from the
27 human genome (hg38) with the default parameters (Frazee et al. 2015). We provided
28 the full-length transcript structure for IPAfinder to infer and quantify IPA sites based
29 on the synthetic RNA-seq dataset.

30

31 **Comparison between IPAfinder-analyzed RNA-seq and custom-analyzed 3'-seq**

32 A total of 330 recurrent CLL-lpa events were obtained from the datasets of Lee et al.
33 (Lee et al. 2018). When IPAfinder was applied to RNA-seq data, an lpa isoform was
34 considered as a recurrently upregulated lpa isoform if it had significant upregulation in
35 at least three malignant B cell samples (11 samples in total) compared with the level in
36 normal immune cell samples. With this criterion, we obtained 306 recurrently
37 upregulated lpa events. A lower Ratio_{MSE} value means that there is a better coverage

1 segmentation point in the given intron region. Thus, CLL samples with a larger
2 number of CLL-IpA events as reported by the original publication (including CLL4,
3 CLL7, CLL11, and CLL12) have more low Ratio_{MSE} values than samples with a
4 smaller number of CLL-IpA events or normal samples (**Supplemental Fig. S4B**),
5 which suggests that Ratio_{MSE} is a rational index for identifying potential intronic poly(A)
6 sites. Furthermore, CLL samples with a larger number of CLL-IpA events have higher
7 IPU (Supplemental Fig. S4B), consistent with previous results (Lee et al. 2018),
8 which indicates that IPU is also a rational index for quantifying IpA isoform usage.

9 10 **Motif frequency analysis**

11 The genomic sequences (from the human reference sequence hg38) of 200
12 nucleotides upstream and downstream of the cryptic 3' splice sites were used for motif
13 analysis. The frequency of HNRNPC binding motif (U)₅ tracts was calculated by
14 counting the number of (U)₅ motifs (smoothed by ± 5 nucleotides centered on the
15 position of interest) along these specified annotation features.

16 17 **Clinical significance analysis of IpA usage**

18 We obtained clinical information including overall survival time of patients from the
19 GDC data portal (<https://portal.gdc.cancer.gov/>). A log-rank test and Kaplan–Meier
20 survival analysis were performed to identify the association between intronic pA site
21 usage and overall survival. For the gene *TSC1*, groups with high and low IpA usage
22 were separately used for a survival plot by splitting the ordered IPU with an equal
23 number of samples in each group. For the gene *CCND2*, patients whose *CCND2* IpA
24 usage is greater than 0.1 were grouped into patients with *CCND2* IpA. All statistical
25 analyses were performed in R (v.3.5.1) (R Core Team 2018).

26 27 **RT-PCR validation of upregulated IpA isoforms upon knockdown of RBPs**

28 Endogenous PTBP1 and HNRNPC were knocked down using pLKO.1-puro lentiviral
29 vector-based shRNAs (Sigma-Aldrich). HEK293T cells were transduced in six-well
30 plates using Lipofectamine 2000 (Invitrogen). Virus was produced using the helper
31 plasmids VSVG and gag/pol. After infection over 36 hours (h), the cells were selected
32 with puromycin (2.5 μ g/ml) for 2 days and the surviving cells were cultured for two
33 more days and then collected for RT–PCR analysis.

34 Total RNA was extracted with TRIzol reagent (Invitrogen) according to the
35 manufacturer's instruction. RNA was reversely transcribed using the FastKing RT Kit
36 (With gDNase) (Tiangen). 20 μ l cDNA product was diluted 5-fold and 2 μ l diluted
37 cDNA was used as the template for each semi-quantitative RT-PCR reaction. We

1 used a typical reaction containing 500 nM forward and reverse primers for individual
2 isoforms. The PCR reaction products were analyzed by gel electrophoresis. Primers
3 are listed in **Supplemental Table S3**.

4

5 **Vector construction**

6 The full-length *TSC1*, *TSC2*, and *SPRED2* mRNA was amplified from HEK293T cDNA.
7 Plasmids for the expression of full length of HA-TSC1 (ENSG00000165699, 1164aa),
8 FLAG-TSC2 (ENSG00000103197, 1807aa) and HA-SPRED2 (ENSG00000198369,
9 418aa) were constructed by cloning full-length CDS of *TSC1*, *TSC2* and *SPRED2* into
10 the pRK5 vector with either FLAG or HA tag at their N terminus. *TSC1* IpA was
11 PCR-amplified from two fragments. Fragment 1 was amplified from HEK293T cDNA
12 and corresponds to amino acids 1–421, whereas fragment 2 was amplified from
13 genomic DNA of HEK293T and corresponds to intronic sequence upstream predicted
14 IPA site.

15 To construct the pCDH-*SPRED2* plasmid, full-length CDS of *SPRED2* cloned from
16 HEK293T cDNA was inserted into the pCDH-CMV-MCS-T2Apuro plasmid using
17 EcoRI/BamHI restriction sites. *SPRED2* IpA was also PCR-amplified from two
18 fragments. Fragment 1 was amplified from HEK293T cDNA and corresponds to amino
19 acids 1–68, whereas fragment 2 was amplified from genomic DNA and corresponds to
20 intronic sequence upstream predicted IPA site. The integrity of all constructs was
21 confirmed by Sanger sequencing.

22

23 **Western blotting**

24 Cells were rinsed with PBS and lysed in cold RIPA buffer (25 mM Tris pH 7.6, 150 mM
25 NaCl, 1% NP-40, 1% sodium deoxycholate, 0.1% SDS) containing freshly added
26 Protease and Phosphatase Inhibitor Cocktail, EDTA-free (Thermo Fisher Scientific).
27 Cell lysates were incubated on ice for 10 minutes, and centrifuged at 14,000 × g for 15
28 minutes at 4 °C. The supernatant was collected and the protein concentration was
29 determined by Bicinchoninic Acid (BCA) assay (Beyotime). A total of 20 µg protein per
30 sample was resolved by 10% SDS-PAGE, followed by transfer to a PVDF membrane
31 with pore size 0.2 µm (Millipore) for immunoblotting. Quantification was performed by
32 densitometry using ImageJ software, and ACTB served as internal control.

33 The following primary antibodies were used: anti-ACTB (proteintech, HRP-60008,
34 1:2000), Anti-phospho-S6 ribosomal protein (Cell signaling, 2215S, 1:2000), Anti-S6
35 ribosomal protein (Cell signaling, 2217S, 1:2000), anti HA-Tag (ABclonal, AE008,
36 1:2000), Mouse anti DDDDK-Tag (ABclonal, AE005, 1:2000). The secondary
37 antibodies used included HRP Goat Anti-Rabbit IgG (H+L) (ABclonal, AS014, 1:5000)

1 and HRP Goat Anti-Mouse IgG (H+L) (ABclonal, AS003, 1:5000).

2

3 **Cell proliferation assay**

4 Cells were counted and seeded in 96-well plates with 2000 cells per well and four
5 replicates for each time points. Cell Counting Kit-8 (CCK-8) reagent (Dojindo) was
6 diluted with medium according to the manufacturer's protocol and then added to each
7 testing wells. Then, cells were incubated at 37 °C for another 2 h and then the
8 absorbance of each well was measured at 450 nm by a microplate reader (Bio-Rad).

9

10 **Data access**

11 The IPAFinder method is freely available at <https://github.com/ZhaozzReal/IPAFinder>
12 and also in Supplemental Code. The *U2AF2*-KD RNA-seq data generated in this
13 study have been submitted to the NCBI BioProject database
14 (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA660570.

15

16 **Competing interest statement**

17 The authors declare no competing interests.

18

19 **Acknowledgments**

20 We thank Edanz Group China (<http://www.liwenbianji.cn/ac>), for editing the English
21 text of a draft of this manuscript. This work was supported by National Key Research
22 and Development Program of China [2018YFC1003500], National Natural Science
23 Foundation of China [91949107, 31771336, 31521003], Shanghai Municipal Science
24 and Technology Major Project (2017SHZDZX01).

25

26 **References**

- 27 Anders S, Reyes A, Huber W. 2012. Detecting differential usage of exons from RNA-seq data. *Genome*
28 *Res* **22**: 2008-2017.
- 29 Attig J, Agostini F, Gooding C, Chakrabarti AM, Singh A, Haberman N, Zagalak JA, Emmett W, Smith CWJ,
30 Luscombe NM et al. 2018. Heteromeric RNP Assembly at LINEs Controls Lineage-Specific RNA
31 Processing. *Cell* **174**: 1067-1081.e1017.
- 32 Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS. 2009. MEME
33 SUITE: tools for motif discovery and searching. *Nucleic acids research* **37**: W202-W208.
- 34 Berg MG, Singh LN, Younis I, Liu Q, Pinto AM, Kaida D, Zhang Z, Cho S, Sherrill-Mix S, Wan L et al. 2012.
35 U1 snRNP determines mRNA length and regulates isoform expression. *Cell* **150**: 53-64.
- 36 Berkovits BD, Mayr C. 2015. Alternative 3' UTRs act as scaffolds to regulate membrane protein
37 localization. *Nature* **522**: 363-367.
- 38 Chen M, Lyu G, Han M, Nie H, Shen T, Chen W, Niu Y, Song Y, Li X, Li H et al. 2018. 3' UTR lengthening as
39 a novel mechanism in regulating cellular senescence. *Genome Res*

- 1 doi:10.1101/gr.224451.117.
- 2 Chong-Kopera H, Inoki K, Li Y, Zhu T, Garcia-Gonzalo FR, Rosa JL, Guan KL. 2006. TSC1 stabilizes TSC2 by
3 inhibiting the interaction between TSC2 and the HERC1 ubiquitin ligase. *J Biol Chem* **281**:
4 8313-8316.
- 5 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. 2013.
6 STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15-21.
- 7 Dubbury SJ, Boutz PL, Sharp PA. 2018. CDK12 regulates DNA repair genes by suppressing intronic
8 polyadenylation. *Nature* **564**: 141-145.
- 9 Elkon R, Ugalde AP, Agami R. 2013. Alternative cleavage and polyadenylation: extent, regulation and
10 function. *Nat Rev Genet* **14**: 496-506.
- 11 Frazee AC, Jaffe AE, Langmead B, Leek JT. 2015. Polyester: simulating RNA-seq datasets with
12 differential transcript expression. *Bioinformatics* **31**: 2778-2784.
- 13 Garami A, Zwartkruis FJ, Nobukuni T, Joaquin M, Rocco M, Stocker H, Kozma SC, Hafen E, Bos JL,
14 Thomas G. 2003. Insulin activation of Rheb, a mediator of mTOR/S6K/4E-BP signaling, is
15 inhibited by TSC1 and 2. *Mol Cell* **11**: 1457-1466.
- 16 Gruber AJ, Schmidt R, Ghosh S, Martin G, Gruber AR, van Nimwegen E, Zavolan M. 2018. Discovery of
17 physiological and cancer-related regulators of 3' UTR processing with KAPAC. *Genome Biol* **19**:
18 44.
- 19 Gruber AJ, Schmidt R, Gruber AR, Martin G, Ghosh S, Belmadani M, Keller W, Zavolan M. 2016. A
20 comprehensive analysis of 3' end sequencing data sets reveals novel polyadenylation signals
21 and the repressive role of heterogeneous ribonucleoprotein C on cleavage and
22 polyadenylation. *Genome research* **26**: 1145-1159.
- 23 Gueroussov S, Gonatopoulos-Pournatzis T, Irimia M, Raj B, Lin ZY, Gingras AC, Blencowe BJ. 2015. An
24 alternative splicing event amplifies evolutionary differences between vertebrates. *Science*
25 **349**: 868-873.
- 26 Guvenc A, Tian B. 2018. Analysis of alternative cleavage and polyadenylation in mature and
27 differentiating neurons using RNA-seq data. *Quant Biol* **6**: 253-266.
- 28 Ha KCH, Blencowe BJ, Morris Q. 2018. QAPA: a new method for the systematic analysis of alternative
29 polyadenylation from RNA-seq data. *Genome Biol* **19**: 45.
- 30 Herrmann CJ, Schmidt R, Kanitz A, Artimo P, Gruber AJ, Zavolan M. 2020. PolyASite 2.0: a consolidated
31 atlas of polyadenylation sites from 3' end sequencing. *Nucleic acids research* **48**: D174-D179.
- 32 Hilgers V, Lemke SB, Levine M. 2012. ELAV mediates 3' UTR extension in the Drosophila nervous
33 system. *Genes Dev* **26**: 2259-2264.
- 34 Hoque M, Ji Z, Zheng D, Luo W, Li W, You B, Park JY, Yehia G, Tian B. 2013. Analysis of alternative
35 cleavage and polyadenylation by 3' region extraction and deep sequencing. *Nat Methods* **10**:
36 133-139.
- 37 Inoki K, Li Y, Xu T, Guan KL. 2003. Rheb GTPase is a direct target of TSC2 GAP activity and regulates
38 mTOR signaling. *Genes Dev* **17**: 1829-1834.
- 39 Kaida D, Berg MG, Younis I, Kasim M, Singh LN, Wan L, Dreyfuss G. 2010. U1 snRNP protects
40 pre-mRNAs from premature cleavage and polyadenylation. *Nature* **468**: 664-668.
- 41 Kasowitz SD, Ma J, Anderson SJ, Leu NA, Xu Y, Gregory BD, Schultz RM, Wang PJ. 2018. Nuclear m6A
42 reader YTHDC1 regulates alternative polyadenylation and splicing during mouse oocyte
43 development. *PLoS Genet* **14**: e1007412.
- 44 Kladney RD, Cardiff RD, Kwiatkowski DJ, Chiang GG, Weber JD, Arbeit JM, Lu ZH. 2010. Tuberous

- 1 sclerosis complex 1: an epithelial tumor suppressor essential to prevent spontaneous
2 prostate cancer in aged mice. *Cancer Res* **70**: 8937-8947.
- 3 König J, Zarnack K, Rot G, Curk T, Kaykci M, Zupan B, Turner DJ, Luscombe NM, Ule J. 2010. iCLIP
4 reveals the function of hnRNP particles in splicing at individual nucleotide resolution. *Nat*
5 *Struct Mol Biol* **17**: 909-915.
- 6 Kyburz A, Friedlein A, Langen H, Keller W. 2006. Direct interactions between subunits of CPSF and the
7 U2 snRNP contribute to the coupling of pre-mRNA 3' end processing and splicing. *Mol Cell* **23**:
8 195-205.
- 9 Lee SH, Singh I, Tisdale S, Abdel-Wahab O, Leslie CS, Mayr C. 2018. Widespread intronic
10 polyadenylation inactivates tumour suppressor genes in leukaemia. *Nature* **561**: 127-131.
- 11 Li W, You B, Hoque M, Zheng D, Luo W, Ji Z, Park JY, Gunderson SI, Kalsotra A, Manley JL et al. 2015.
12 Systematic profiling of poly(A)+ transcripts modulated by core 3' end processing and splicing
13 factors reveals regulatory rules of alternative cleavage and polyadenylation. *PLoS Genet* **11**:
14 e1005166.
- 15 Licatalosi DD, Mele A, Fak JJ, Ule J, Kaykci M, Chi SW, Clark TA, Schweitzer AC, Blume JE, Wang X et al.
16 2008. HITS-CLIP yields genome-wide insights into brain alternative RNA processing. *Nature*
17 **456**: 464-469.
- 18 Liu N, Dai Q, Zheng G, He C, Parisien M, Pan T. 2015. N(6)-methyladenosine-dependent RNA structural
19 switches regulate RNA-protein interactions. *Nature* **518**: 560-564.
- 20 Mayr C, Bartel DP. 2009. Widespread shortening of 3'UTRs by alternative cleavage and
21 polyadenylation activates oncogenes in cancer cells. *Cell* **138**: 673-684.
- 22 Millevoi S, Loulergue C, Dettwiler S, Karaa SZ, Keller W, Antoniou M, Vagner S. 2006. An interaction
23 between U2AF 65 and CF I(m) links the splicing and 3' end processing machineries. *Embo j* **25**:
24 4854-4864.
- 25 Mirzaa G, Parry DA, Fry AE, Giamanco KA, Schwartzenruber J, Vanstone M, Logan CV, Roberts N,
26 Johnson CA, Singh S et al. 2014. De novo CCND2 mutations leading to stabilization of cyclin
27 D2 cause megalencephaly-polymicrogyria-polydactyly-hydrocephalus syndrome. *Nat Genet*
28 **46**: 510-515.
- 29 Mueller AA, van Velthoven CT, Fukumoto KD, Cheung TH, Rando TA. 2016. Intronic polyadenylation of
30 PDGFR α in resident stem cells attenuates muscle fibrosis. *Nature* **540**: 276-279.
- 31 Ni T, Yang Y, Hafez D, Yang W, Kiesewetter K, Wakabayashi Y, Ohler U, Peng W, Zhu J. 2013. Distinct
32 polyadenylation landscapes of diverse human tissues revealed by a modified PA-seq strategy.
33 *BMC Genomics* **14**: 615.
- 34 Nobuhisa I, Kato R, Inoue H, Takizawa M, Okita K, Yoshimura A, Taga T. 2004. Spred-2 suppresses
35 aorta-gonad-mesonephros hematopoiesis by inhibiting MAP kinase activation. *J Exp Med* **199**:
36 737-742.
- 37 Oberstrass FC, Auweter SD, Erat M, Hargous Y, Henning A, Wenter P, Reymond L, Amir-Ahmady B,
38 Pitsch S, Black DL et al. 2005. Structure of PTB bound to RNA: specific binding and
39 implications for splicing regulation. *Science* **309**: 2054-2057.
- 40 Oh JM, Di C, Venters CC, Guo J, Arai C, So BR, Pinto AM, Zhang Z, Wan L, Younis I et al. 2017. U1 snRNP
41 telescripting regulates a size-function-stratified human genome. *Nat Struct Mol Biol* **24**:
42 993-999.
- 43 R Core Team. 2018. R: a language and environment for statistical computing. R Foundation for
44 Statistical Computing, Vienna. <https://www.R-project.org/>.

- 1 Shepard PJ, Choi E-A, Lu J, Flanagan LA, Hertel KJ, Shi Y. 2011. Complex and dynamic landscape of RNA
2 polyadenylation revealed by PAS-Seq. *RNA (New York, NY)* **17**: 761-772.
- 3 Singh I, Lee SH, Sperling AS, Samur MK, Tai YT, Fulcinitti M, Munshi NC, Mayr C, Leslie CS. 2018.
4 Widespread intronic polyadenylation diversifies immune cell transcriptomes. *Nat Commun* **9**:
5 1716.
- 6 Śledź P, Jinek M. 2016. Structural insights into the molecular mechanism of the m(6)A writer complex.
7 *Elife* **5**.
- 8 So BR, Di C, Cai Z, Venters CC, Guo J, Oh JM, Arai C, Dreyfuss G. 2019. A Complex of U1 snRNP with
9 Cleavage and Polyadenylation Factors Controls Telescripting, Regulating mRNA Transcription
10 in Human Cells. *Mol Cell* **76**: 590-599.e594.
- 11 Sun S, Chen S, Liu F, Wu H, McHugh J, Bergin IL, Gupta A, Adams D, Guan JL. 2015. Constitutive
12 Activation of mTORC1 in Endothelial Cells Leads to the Development and Progression of
13 Lymphangiosarcoma through VEGF Autocrine Signaling. *Cancer Cell* **28**: 758-772.
- 14 Takano Y, Kato Y, Masuda M, Ohshima Y, Okayasu I. 1999. Cyclin D2, but not cyclin D1, overexpression
15 closely correlates with gastric cancer progression and prognosis. *J Pathol* **189**: 194-200.
- 16 Takano Y, Kato Y, van Diest PJ, Masuda M, Mitomi H, Okayasu I. 2000. Cyclin D2 overexpression and
17 lack of p27 correlate positively and cyclin E inversely with a poor prognosis in gastric cancer
18 cases. *Am J Pathol* **156**: 585-594.
- 19 Tian B, Manley JL. 2017. Alternative polyadenylation of mRNA precursors. *Nat Rev Mol Cell Biol* **18**:
20 18-30.
- 21 Tian B, Pan Z, Lee JY. 2007. Widespread mRNA polyadenylation events in introns indicate dynamic
22 interplay between polyadenylation and splicing. *Genome Res* **17**: 156-165.
- 23 van Slegtenhorst M, Nellist M, Nagelkerken B, Cheadle J, Snell R, van den Ouweland A, Reuser A,
24 Sampson J, Halley D, van der Sluijs P. 1998. Interaction between hamartin and tuberlin, the
25 TSC1 and TSC2 gene products. *Hum Mol Genet* **7**: 1053-1057.
- 26 Wakioka T, Sasaki A, Kato R, Shouda T, Matsumoto A, Miyoshi K, Tsuneoka M, Komiya S, Baron R,
27 Yoshimura A. 2001. Spred is a Sprouty-related suppressor of Ras signalling. *Nature* **412**:
28 647-651.
- 29 Wang R, Nambiar R, Zheng D, Tian B. 2018. PolyA_DB 3 catalogs cleavage and polyadenylation sites
30 identified by deep sequencing in multiple genomes. *Nucleic Acids Res* **46**: D315-d319.
- 31 Wang R, Tian B. 2020. APALyzer: a bioinformatics package for analysis of alternative polyadenylation
32 isoforms. *Bioinformatics* **36**: 3907-3909.
- 33 Wu S, Romfo CM, Nilsen TW, Green MR. 1999. Functional recognition of the 3' splice site AG by the
34 splicing factor U2AF35. *Nature* **402**: 832-835.
- 35 Xia Z, Donehower LA, Cooper TA, Neilson JR, Wheeler DA, Wagner EJ, Li W. 2014. Dynamic analyses of
36 alternative polyadenylation from RNA-seq reveal a 3'-UTR landscape across seven tumour
37 types. *Nat Commun* **5**: 5274.
- 38 Xiao W, Adhikari S, Dahal U, Chen YS, Hao YJ, Sun BF, Sun HY, Li A, Ping XL, Lai WY et al. 2016. Nuclear
39 m(6)A Reader YTHDC1 Regulates mRNA Splicing. *Mol Cell* **61**: 507-519.
- 40 Yang H, Yu Z, Chen X, Li J, Li N, Cheng J, Gao N, Yuan HX, Ye D, Guan KL et al. 2021. Structural insights
41 into TSC complex assembly and GAP activity on Rheb. *Nat Commun* **12**: 339.
- 42 Yang SW, Li L, Connelly JP, Porter SN, Kodali K, Gan H, Park JM, Tacer KF, Tillman H, Peng J et al. 2020. A
43 Cancer-Specific Ubiquitin Ligase Drives mRNA Alternative Polyadenylation by Ubiquitinating
44 the mRNA 3' End Processing Complex. *Mol Cell* **77**: 1206-1221.e1207.

- 1 Yang Y, Hsu PJ, Chen YS, Yang YG. 2018. Dynamic transcriptomic m(6)A decoration: writers, erasers,
2 readers and functions in RNA metabolism. *Cell Res* **28**: 616-624.
- 3 Yao J, Ding D, Li X, Shen T, Fu H, Zhong H, Wei G, Ni T. 2020. Prevalent intron retention fine-tunes gene
4 expression and contributes to cellular senescence. *Aging Cell* **19**: e13276.
- 5 Ye C, Long Y, Ji G, Li QQ, Wu X. 2018. APAtap: identification and quantification of alternative
6 polyadenylation sites from RNA-seq data. *Bioinformatics* **34**: 1841-1849.
- 7 Zamore PD, Patton JG, Green MR. 1992. Cloning and domain structure of the mammalian splicing
8 factor U2AF. *Nature* **355**: 609-614.
- 9 Zarnack K, Konig J, Tajnik M, Martincorena I, Eustermann S, Stevant I, Reyes A, Anders S, Luscombe
10 NM, Ule J. 2013. Direct competition between hnRNP C and U2AF65 protects the
11 transcriptome from the exonization of Alu elements. *Cell* **152**: 453-466.

12
13
14
15

16 **Figure 1.** Overview of the IPAFinder algorithm and evaluation of its performance. (A)
17 Schematic diagram of IPAFinder in detecting composite terminal exon lpa event.
18 Intronic poly(A) site is determined based on the expected drop in read coverage
19 downstream the predicted poly(A) site. Alternative splicing events are excluded by
20 recognizing junction-spanning reads. (B,C) Two examples of IPAFinder-identified
21 dynamically changed lpa events from TCGA RNA-seq data. lpa usage of the *ELP5*
22 gene (B, composite lpa site) and *ANOS1* gene (C, skipped lpa site) is increased in
23 both lung squamous cell carcinoma (LUSC) and lung adenocarcinoma (LUAD)
24 compared with that in matched normal tissues. Sample IDs are shown at the top-right
25 corner of the corresponding RNA-seq density plot. lpa site is indicated by a red arrow.
26 (D) Performance of IPAFinder in detecting differentially used lpa sites in terms of
27 sensitivity and precision. The number of TPs, FPs and predefined true differentially
28 used lpa sites (P) are used to calculate sensitivity (TP/P) and precision (TP/(TP+FP)).
29 For each coverage level, we repeated ten times to calculate the mean value of
30 sensitivity and precision. For samples without replicates, two methods including
31 bootstrapping-based method and Fisher's exact test-based were assessed. (E) Venn
32 diagram comparison of the number of recurrent upregulated lpa events identified by
33 IPAFinder and those by 3'-seq using the same data from CLL and immune cell
34 samples. (F) An example of dynamically changed lpa events (*PDXDC1*) between CLL
35 samples and normal B cells detected by IPAFinder, which was absent in 3'-seq.

36

37 **Figure 2.** IPAFinder reveals the global landscape of lpa events across six TCGA
38 tumor types. (A) IPAFinder discovers prevalent upregulation of lpa events across six
39 tumor types. The upper histogram shows the number of dynamically changed lpa

1 events per tumor. The lower heatmap shows IpA events (rows) undergoing
 2 upregulation (red) or downregulation (blue) in each of the 256 tumors (columns)
 3 compared with the levels in matched normal tissues across six tumor types. (B) Bar
 4 plots showing the percentages of IPAFinder-predicted breakpoints (left) and the
 5 randomly selected positions (right) that overlap with annotated intronic pA sites
 6 (RefSeq, UCSC, Ensembl, PolyASite 2.0). The *P*-value was calculated by *t*-test using
 7 100× bootstrapping of data. (C) MEME identifies the enrichment of the canonical
 8 poly(A) signal AATAAA in the ± 50 bp region around IPAFinder-inferred IPA sites. The
 9 corresponding genomic sequences (from human reference sequence hg38) serve as
 10 input. (D) The distribution of the retained coding region fraction (resulting from IpA
 11 usage) of the annotated longest coding region (CDR). (E) Box plot for lengths of
 12 introns with skipped IpA, composite IpA, and introns without IpA. The *P*-value was
 13 calculated based on two-sided Wilcoxon rank-sum test. (F) *AUH* as an example to
 14 display skipped IPA sites with increased usage in two types of cancer (LUAD and
 15 HNSC) located in an extremely long intron. Sample IDs are shown at the top-right
 16 corner of the corresponding RNA-seq density plot.

17
 18

19 **Figure 3.** IpA generates truncated proteins with important functional impacts. (A)
 20 RNA-seq density plots showing that *TSC1* has increased IpA usage in lung cancers.
 21 Numbers on the y-axis indicate RNA-seq read coverage. Sample IDs are shown at
 22 the top-right corner of the corresponding RNA-seq density plot. (B) Box plots showing
 23 that *TSC1* has significantly higher IpA usage in LUAD, LUSC, and HNSC tumors. The
 24 *P*-value was calculated based on two-sided Wilcoxon rank-sum test. (C) Immunoblot
 25 analysis of S6 phosphorylation in HEK293T cells with overexpression of the IpA and
 26 full-length (FL) isoform of *TSC1*. Successful expression of HA-tagged IpA and
 27 full-length isoforms of *TSC1* was confirmed by Western blot using anti-HA and
 28 anti-FLAG antibodies (left). Both total S6 protein and its phosphorylated form (P-S6)
 29 were also quantified by western blot (right). ACTB was used as an internal control. NC
 30 (negative control) means a sample without FLAG-TSC2-WT, HA-TSC1-FL and
 31 HA-TSC1-IpA. *** denotes *P*-value < 0.001, *t*-test. (D) Kaplan–Meier curves of overall
 32 survival for two HNSC tumor groups stratified by the IpA usage of *TSC1*. The *P*-value
 33 was calculated using the log-rank test. (E) Diagrams showing the domain information
 34 of full-length and IpA-generated truncated proteins, with known domains shown in
 35 green. The numbers of retained and novel amino acids (aa) and amino acids of
 36 full-length proteins are given. The position of the important residue Thr280 of CCND2
 37 is denoted by a short blue line (p.Thr280). (F) RNA-seq density plots showing that

1 *SPRED2* has increased lpa usage in lung cancers. (G) CCK-8 assay showing that
 2 lpa-truncated *SPRED2* fails to inhibit cell proliferation in NCI-H520 cells. *** denotes
 3 *P*-value < 0.001, *t*-test. (H) RNA-seq density plots showing that *CCND2* has increased
 4 lpa usage in lung cancers.

5

6 **Figure 4.** PTBP1/2 and HNRNPC inhibit the usage of IPA sites. (A) Scatterplot of IPUI
 7 value reflecting the relative lpa usage before and after concurrent knockdown of
 8 *PTBP1* and *PTBP2* (*PTBP1/2*-KD) in HEK293 cells. Red and blue dots represent
 9 genes with increased and decreased lpa usage upon knockdown of *PTBP1/2*,
 10 respectively. (B) Representative examples of lpa events repressed by PTBP1/2. Both
 11 skipped lpa event (top) and composite lpa event (bottom) are shown. (C) Sequences
 12 flanking example IPA sites repressed by PTBP1/2 have CU repeats (red) and poly(A)
 13 signals (blue). Five genes with a skipped terminal exon (top) and five genes with a
 14 composite terminal exon (bottom) are shown. The first bases of skipped terminal
 15 exons are denoted by enlarged characters. (D) Scatterplot of IPUI value reflecting the
 16 relative lpa usage before and after knockdown of *HNRNPC* (*HNRNPC*-KD) in
 17 HEK293T cells. Red and blue dots represent genes with increased and decreased lpa
 18 usage in *HNRNPC*-KD cells, respectively. (E) Representative examples of
 19 HNRNPC-repressed lpa events. (F) Sequences flanking example IPA sites repressed
 20 by HNRNPC have (U)₅ tracts (red) and poly(A) signals (blue).

21

22

23 **Figure 5.** IPAFinder reveals that intronic polyadenylation can be influenced by factors
 24 related to splicing and m⁶A modification. (A,C,E,G) Scatterplot of IPUI values
 25 reflecting the relative lpa usage in cells with knockdown (KD) of *U2AF1* (A, in HFF
 26 cell), *U2AF2* (C, in HFF cell), *YTHDC1* (E, in HeLa cell), or *METTL3* (G, in HEK293T
 27 cell). Red and blue dots represent genes with increased and decreased lpa usage
 28 upon knockdown of corresponding genes. (B,D,F,H) Representative RNA-seq density
 29 plots for genes with significantly increased lpa usage upon knockdown of *U2AF1* (B),
 30 *U2AF2* (D), *YTHDC1* (F), or *METTL3* (H). Each knockdown condition has two lpa
 31 examples: the composite terminal exon lpa (top) and the skipped terminal exon lpa
 32 (bottom).

33

34 **Figure 6.** Comparison between IPAFinder and APalyzer. (A) Venn diagram
 35 illustrating the overlap of upregulated lpa events upon *HNRNPC*-KD identified by
 36 IPAFinder and APalyzer, respectively. (B) Representative RNA-seq tracks for genes
 37 with increased lpa usage inferred by IPAFinder but not APalyzer. (C) Sequence

- 1 flanking three IPA sites shown in panel B have (U)₅ tract (red) and poly(A) signal
- 2 (blue). The first bases of skipped terminal exons are denoted by enlarged characters.

