

## Insights into the design and interpretation of iCLIP experiments

Haberman N<sup>1,2</sup>, Huppertz I<sup>1,3,4</sup>, Attig J<sup>1,2</sup>, König J<sup>1,5</sup>, Wang Z<sup>6,7</sup>, Hauer C<sup>4,8,9</sup>, Hentze MW<sup>4,8</sup>, Kulozik AE<sup>8,9</sup>, Le Hir H<sup>6,7</sup>, Curk T<sup>10</sup>, Sibley CR<sup>1,11</sup>, Zarnack K<sup>12</sup>, Ule J<sup>1,2</sup>

<sup>1</sup>Department of Molecular Neuroscience, UCL Institute of Neurology, Queen Square, London, WC1N 3BG, UK

<sup>2</sup>The Crick Institute, 1 Midland Road, London NW1 1AT, UK

<sup>3</sup>MRC Laboratory of Molecular Biology, Francis Crick Avenue, Cambridge, CB2 0QH, UK

<sup>4</sup>European Molecular Biology Laboratory (EMBL), Meyerhofstrasse 1, 69117 Heidelberg, Germany

<sup>5</sup>Institute of Molecular Biology (IMB), Ackermannweg 4, 55128 Mainz, Germany

<sup>6</sup>Institut de Biologie de l'ENS (IBENS), 46 rue d'Ulm, Paris F-75005, France

<sup>7</sup>CNRS UMR 8197, Paris Cedex 05, 75230, France

<sup>8</sup>Molecular Medicine Partnership Unit (MMPU), Im Neuenheimer Feld 350, 69120 Heidelberg, Germany

<sup>9</sup>Department of Pediatric Oncology, Hematology and Immunology, University of Heidelberg, Im Neuenheimer Feld 430, 69120 Heidelberg, Germany

<sup>10</sup>Faculty of Computer and Information Science, University of Ljubljana, Tržaška cesta 25, SI-1000, Ljubljana, Slovenia

<sup>11</sup>Division of Brain Sciences, Department of Medicine, Imperial College London, London, UK.

<sup>12</sup>Buchmann Institute for Molecular Life Sciences (BMLS), Goethe University Frankfurt, Max-von-Laue-Str. 15, 60438 Frankfurt am Main, Germany

### Author Information:

Nejc Haberman and Ina Huppertz contributed equally to this work.

### Corresponding authors:

Jernej Ule: [j.ule@ucl.ac.uk](mailto:j.ule@ucl.ac.uk)

Kathi Zarnack: [kathi.zarnack@bmls.de](mailto:kathi.zarnack@bmls.de)

**Keywords:** protein-RNA interactions, iCLIP, eCLIP, irCLIP, binding site assignment, high-throughput sequencing, PTBP1, eIF4A3, exon-junction complex

## **Additional file 1: Figures S1 to S6**

### **Figure S1 | Quality control of the PTBP1 and eIF4A3 iCLIP.**

- a)** SDS-PAGE analysis of PTBP1-iCLIP2 performed with high RNase conditions. X-ray film shows exposure to the nitrocellulose membrane harbouring PTBP1-RNA complexes that were separated on SDS-PAGE. Size range of the isolated complexes is indicated on the left. MW, molecular weight.
- b)** Same as a), but performed with low RNase conditions.
- c)** Same as a), but showing eIF4A3-RNA complexes isolated in the eIF4A3-iCLIP2 experiment.
- d)** Native acrylamide gel showing eIF4A3-iCLIP2 PCR products that were amplified from medium (M) and high (H) cDNA size ranges. During cDNA library preparation, two size ranges of cDNAs were isolated from the denaturing acrylamide gel as in the recommended protocol [5]. After amplification, 170-270 nt PCR products were obtained, including 128 nt of Illumina primer and iCLIP barcodes, corresponding to an approximate size range of mappable sequenced cDNAs of 40-140 nt.
- d)** Same as a), but this time showing eIF4A3-RNA complexes isolated in the eIF4A3-iCLIP3 experiment.
- e)** Same as d), but this time showing eIF4A3-iCLIP3 PCR products.

### **Figure S2 | CL-motifs are enriched at cDNA deletions and cDNA-starts in U2AF2-iCLIP.**

- a)** Analysis of all U2AF2 experiments examined in this study shows the proportion of cDNAs from each experiment that overlap with a CL-motif at each position relative to the cDNA-start.
- b)** Analysis of all eIF4A3 experiments examined in this study shows the proportion of cDNAs from each experiment that overlap with a CL-motif at each position relative to the cDNA-start.
- c)** Proportion of PTBP1-iCLIP1 cDNAs that overlap with a CL-motif at each position relative to the cDNA-start. Only cDNAs shorter than 40 nt are examined, and are divided into those lacking deletions, or containing a deletion within the first 7 nt, or anywhere in the remaining portion of the cDNA.
- d)** Same as c), but for U2AF2-iCLIP.

### **Figure S3 | Analysis of cDNA-starts and cDNA-ends at the start of Y-tracts.**

- a)** The cDNA-starts of PTBP1-iCLIP1 and CLIP experiments are plotted around the starts of >35 nt Y-tracts that are annotated as T-rich or TC-rich low-complexity sequence in the human genome (hg19). cDNAs of PTBP1-iCLIP1 are divided into four length categories: 17-29 nt, 30-34 nt, 35-39 nt, and >39 nt.
- b)** Same as a), but using U2AF2-iCLIP and CLIP cDNAs.

- c) Same as a), but showing the positions of cDNA-ends.
- d) Same as b), but showing the positions of cDNA-ends.
- e) Same as a), but for the cDNA-centres around the ends of Y-tracks.
- f) Ratio of cDNA-starts and cDNA-centres that are inside of Y-tracks compared to the downstream region (schematic description at the bottom). Statistical test for cDNA-starts and cDNA-centres enrichment in Y-tracks region was done by Fisher's Exact Test with p-value < 2.2e-16.

**Figure S4 | Constrained cDNA-ends in eIF4A3 iCLIP.**

- a) Heatmap showing the position of cDNA-ends around top 1,000 exon-exon junctions contain the highest number of cDNAs in eIF4A3-iCLIP1. Junctions are grouped in each row based on their shared position of cDNA-end peaks (marked by blue rectangle). Each row shows the average of cDNA counts across the grouped junctions. The values are normalised against the maximum value across all rows.
- b) Same as a), but for eIF4A3-iCLIP2.
- c) Same as a), but showing summarised pairing probability around exon-exon junctions in eIF4A3-iCLIP1. The position of the cDNA-end peak is marked by the white rectangle.
- d) Same as c), but for eIF4A3-iCLIP2.
- e) Genomic nucleotide composition around cDNA-end peaks in eIF4A3-iCLIP1.
- f) Same as e), but for eIF4A3-iCLIP2.
- g) Same as e), but for eIF4A3-iCLIP3.
- h) Same as e), but for eIF4A3-CLIP.
- i-k) Distribution of cDNA-starts (solid lines) and ends (dotted lines) in eIF4A3-iCLIP2 relative to exon-exon junctions. Junctions were divided into three different classes according to the position cDNA-end peaks at: -7 to 2 nt (**i**), 3 to 12 nt (**j**), or 13 to 25 nt (**k**) relative to exon-exon junctions. cDNA length categories and labelling as shown on top.

**Figure S5 | The impact of cDNA-end constraints on cDNA-starts in eIF4A3 iCLIP.**

- a) The distribution of cDNA-starts (solid lines) and ends (dotted lines) relative to the cDNA-end peaks that were identified at top 1,000 exon-exon junctions contain the highest number of cDNAs in eIF4A3-iCLIP1. cDNAs are divided into three length categories: 17-29 nt, 30-34 nt, and 35-39 nt.
- b) Same as a), but for eIF4A3-iCLIP3.

- c) Same as a), but for the junction that ranks 8<sup>th</sup> by the number of cDNAs in eIF4A3-iCLIP1.
- d) Same as b), but for the junction that ranks 1<sup>st</sup> by the number of cDNAs in eIF4A3-iCLIP3.
- e) Same as a), but for the junction that ranks 14<sup>th</sup> by the number of cDNAs in eIF4A3-iCLIP1.
- f) Same as b), but for the junction that ranks 4<sup>th</sup> by the number of cDNAs in eIF4A3-iCLIP3.
- g) Same as a), but for the junction that ranks 10<sup>th</sup> by the number of cDNAs in eIF4A3-iCLIP1.
- h) Same as b), but for the junction that ranks 5<sup>th</sup> by the number of cDNAs in eIF4A3-iCLIP3.

**Figure S6 | Distribution of cDNA sizes in the studied experiments.**

- a) Distribution of cDNA sizes in eIF4A3 CLIP and iCLIP experiments of cDNAs that are shorter than 39 nt. The number above the lines reports the % of cDNAs shorter than 39 nt. For longer cDNAs, it is not possible to draw the distribution as their precise lengths are unknown due to the limited length of sequencing. Thus, both the distribution and the % needs to be taken into account to estimate if there is a narrow distribution of cDNA sizes. For example, the distribution shows preferred lengths for both eIF4A3-iCLIP1 and eIF4A3-iCLIP3, but in case of eIF4A3-iCLIP3 only 36% of cDNAs are shorter than 39 nt, while in eIF4A3-iCLIP1 approximately 50% of cDNAs are in the length range of 27-37 nt. Thus, only eIF4A3-iCLIP1 has a strong potential for the cDNA distribution to affect binding site assignment.
- b) Same as a), but for PTBP1 CLIP experiments, showing the % of cDNAs shorter than 34 nt due to the shorter sequencing length.
- c) Same as b), but for U2AF2-iCLIP which shows a trend for shorter cDNA size distribution, with 54% of cDNAs <39 nt. However, this is not a major problem due to the lesser cDNA-end constraints in this experiment (see the broad distribution of cDNA-ends around intron-exon junctions in Figure 6a).

Figure S1

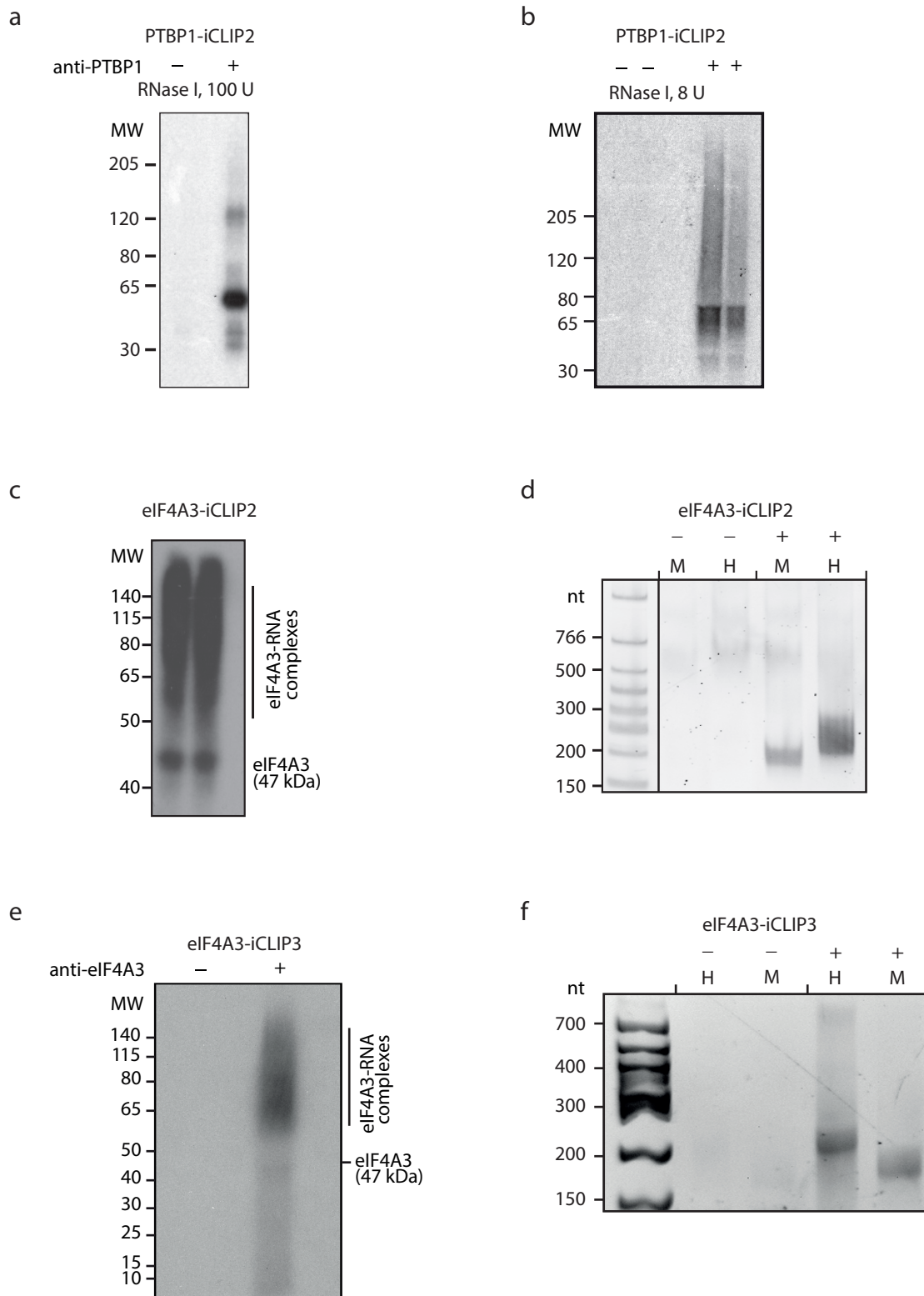
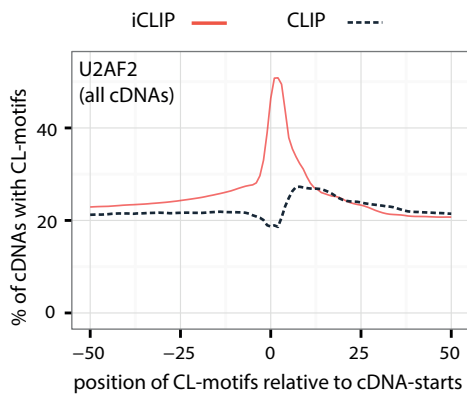
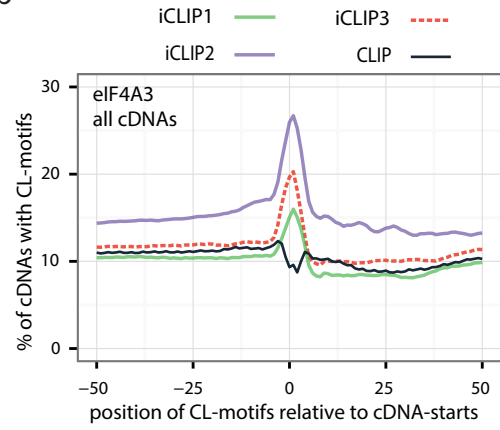


Figure S2

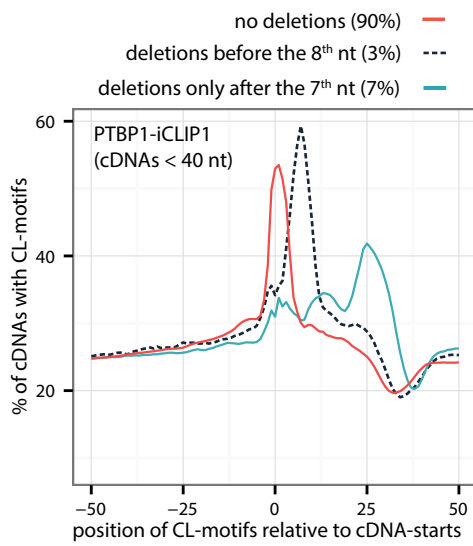
a



b



c



d

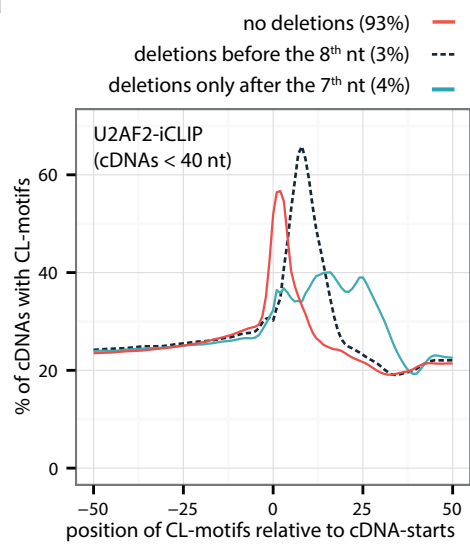


Figure S3

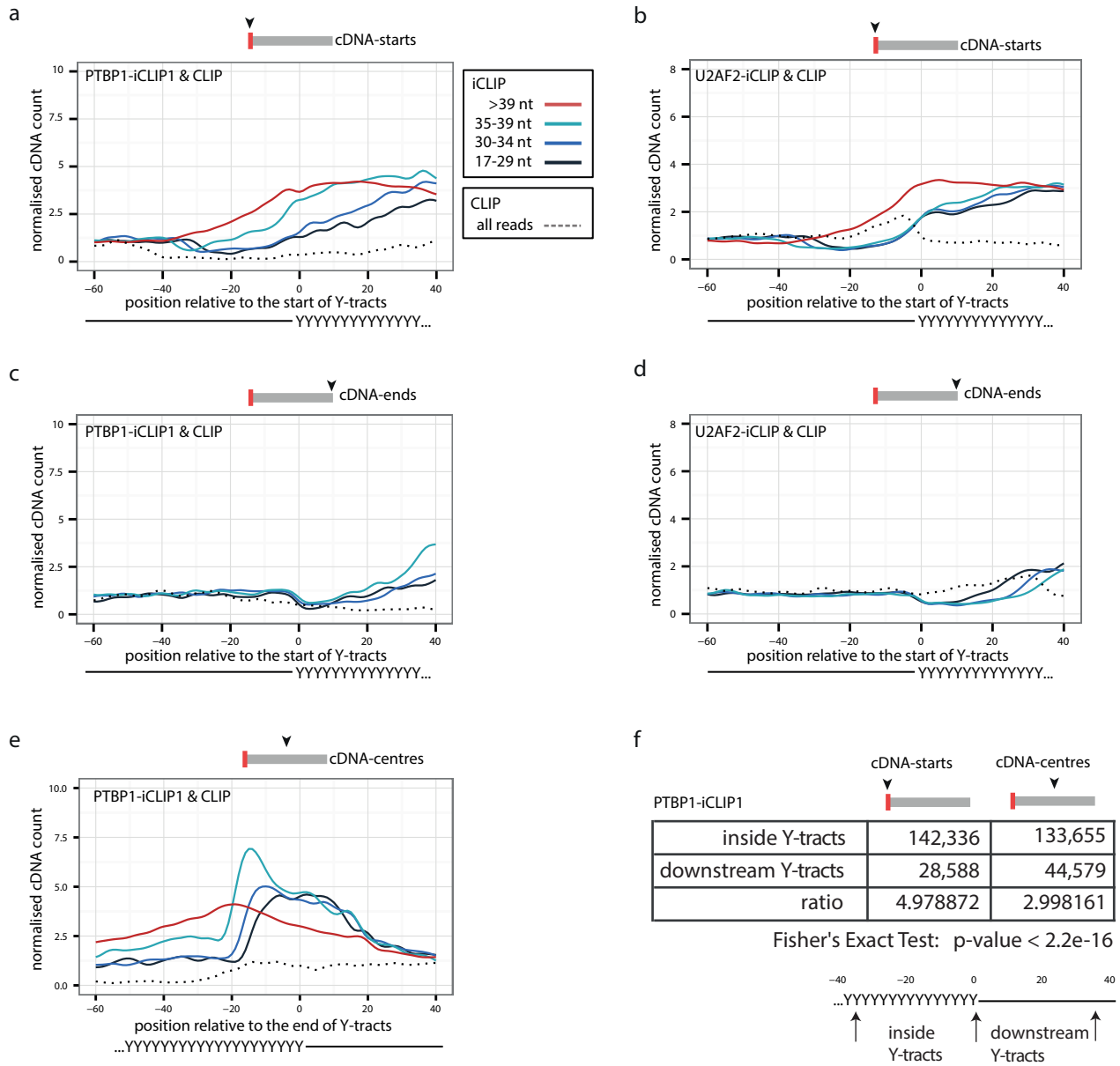


Figure S4

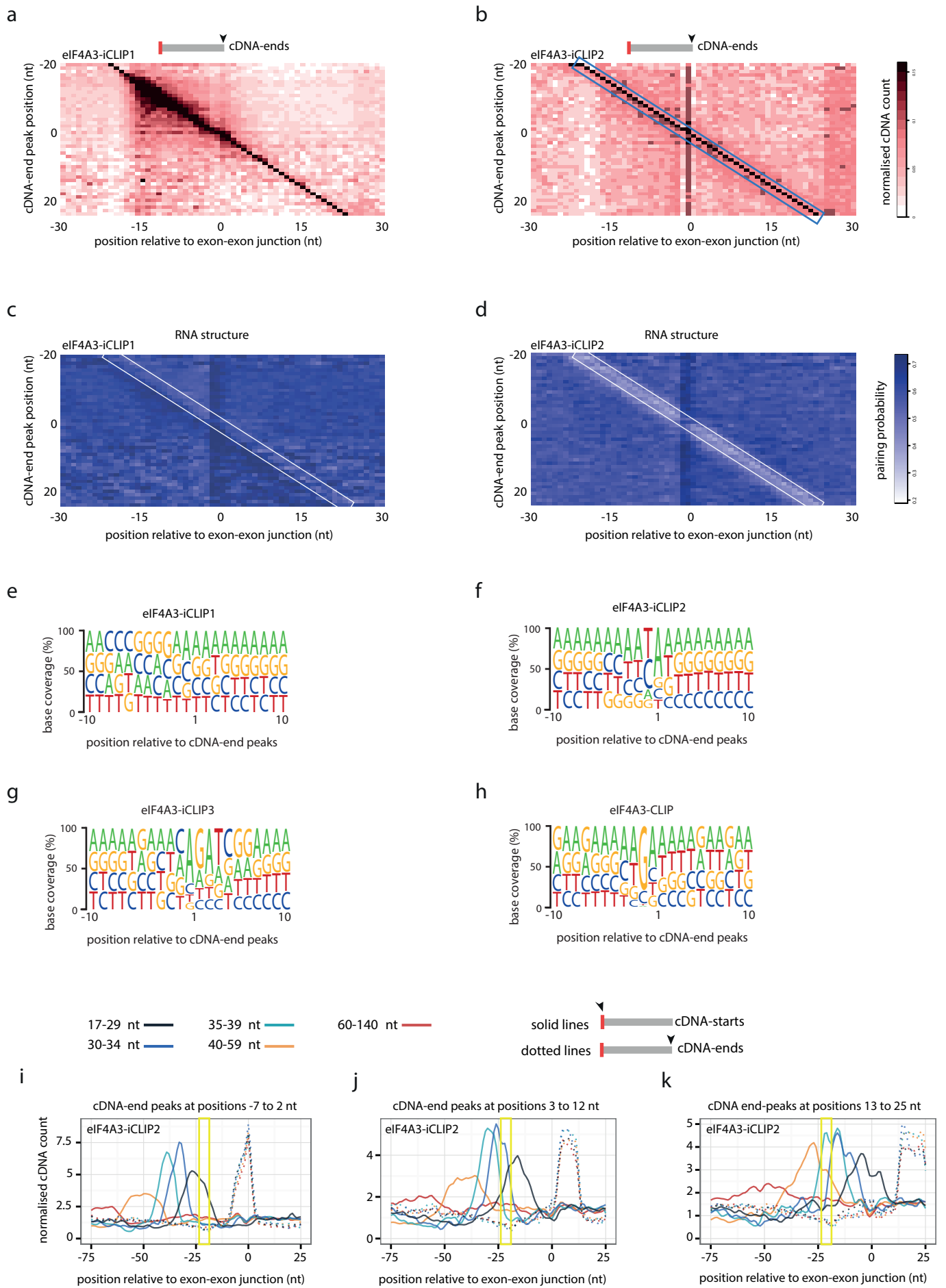




Figure S5

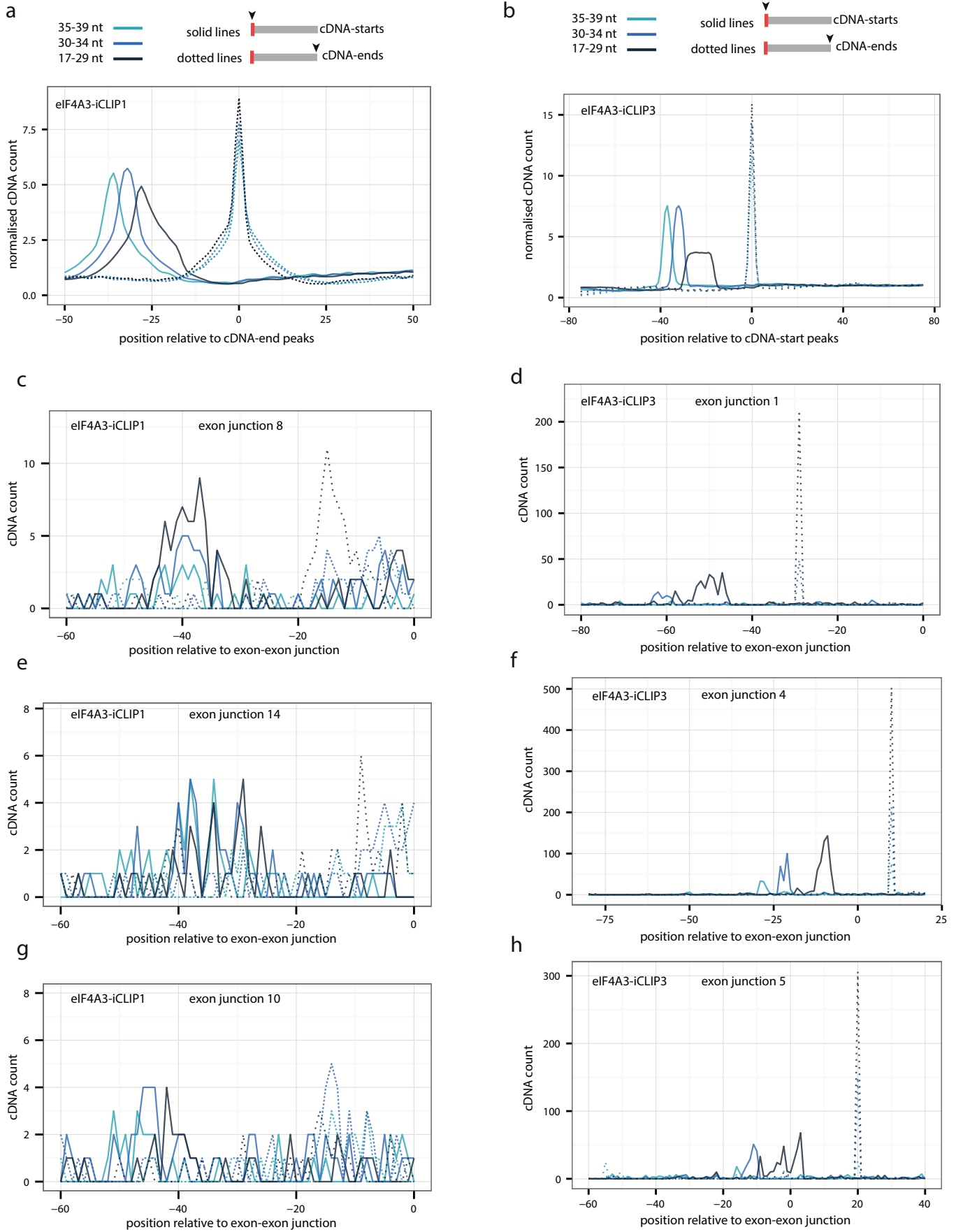
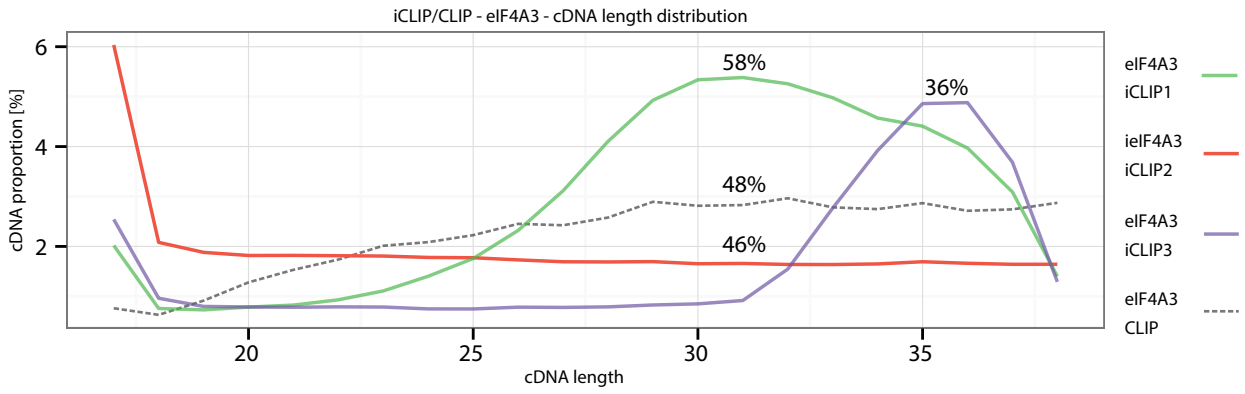
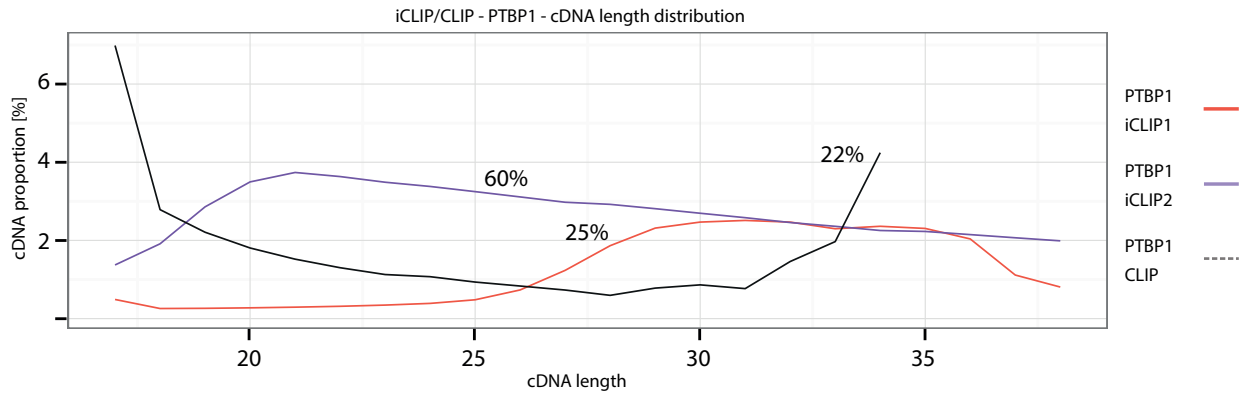


Figure S6

a



b



c

