

Supplementary Data

Preferential formation of Z-RNA over intercalated motifs in long non-coding RNA

Uditi Bhatt¹, Anne Cucchiaroni², Yu Luo², Cameron W. Evans¹, Jean-Louis Mergny², K. Swaminathan Iyer¹, and Nicole M. Smith^{1,*}

¹ School of Molecular Sciences, The University of Western Australia, Crawley, WA, 6009, Australia.

² Laboratoire d'Optique et Biosciences, École Polytechnique, CNRS, INSERM, Institut Polytechnique de Paris, 91128 Palaiseau, France.

* To whom correspondence should be addressed. Tel: +61 8 6488 4423; Email: nicole.smith@uwa.edu.au

Table of Contents

Figure S1. Unfiltered data of predicted non-canonical secondary structures in lncRNA and mRNA.

Table S1. Total number and percentage of predicted structures from lncRNAs and mRNAs with near-equivalent transcript lengths and GC content.

Figure S2. Structure density and percent of transcript data for near-equivalent length and GC% lncRNA and mRNA dataset.

Figure S3. Sequence overlaps between predicted R-loops and predicted G4s and iMs.

Figure S4. Score and tissue expression data for pG4s and piMs in lncRNA.

Figure S5. lncRNA tissue expression.

Figure S6. Score distribution in different tissues.

Figure S7. FRET-MC control data.

Figure S8. Raw IDS data for SNHG14 and MIAT pG4 lncRNA sequences at 37 °C.

Figure S9. UV melt curve raw data and table showing T_m for each pG4 lncRNA oligo tested.

Figure S10. Polyacrylamide gel for pG4 lncRNA.

Figure S11. CD of lncRNA piM (Z-RNA) sequences.

Table S2. pG4 sequences selected for biophysical analysis that are present in multiple lncRNAs

Table S3. Occurrence of lncRNA pG4s within repeat elements of the genome.

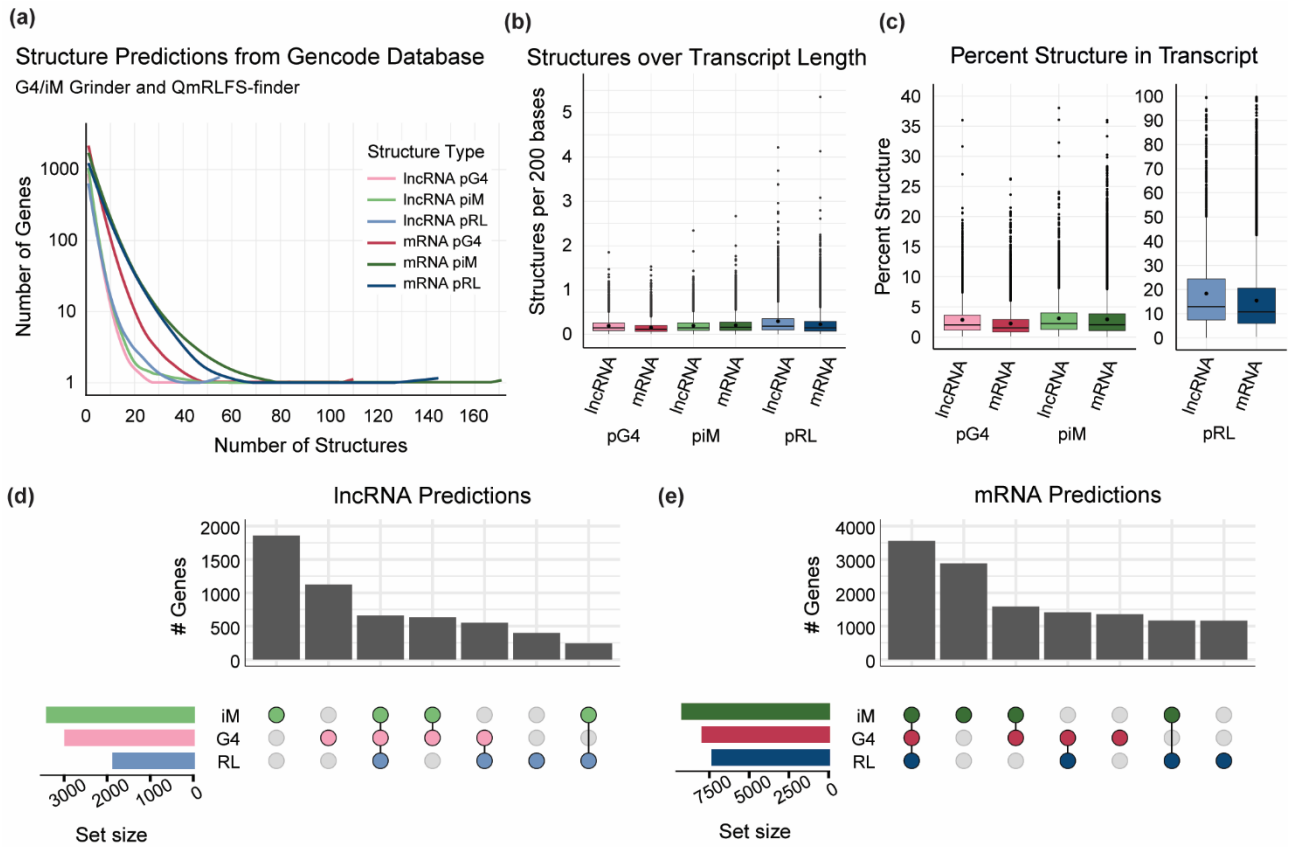


Figure S1. Unfiltered data of predicted non-canonical secondary structures in lncRNA and mRNA. (a) Density of each predicted structure type in lncRNA and mRNA genes. (b) Box plot showing the distribution of the number of structures over transcript length for each type of predicted structure in lncRNA and mRNA with outliers. (c) Box plot showing the distribution of the total percentage of nucleotides in the transcript involved in structure formation for each structure type in lncRNA and mRNA including outliers. (d) lncRNA and (e) mRNA structure prediction datasets. Total number of each structure type shown in horizontal bar graph ‘Set size’.

Table S1. Total number and percentage of pG4, piM, and pRL structures from lncRNA and mRNA with near-equivalent transcript lengths and GC content. 14,637 transcripts from 5,254 genes were included for lncRNA and 13,695 transcripts from 6,260 genes were included for mRNA.

	Structure type	Total predicted structures	Unique genes	% of input genes	Unique transcripts	% of input transcripts	Unique motifs
lncRNA	G4	805	528	10.05	767	5.24	584
	iM	1221	646	12.30	1151	7.86	703
	RL	240	160	3.05	192	1.31	212
mRNA	G4	825	537	8.58	776	5.67	601
	iM	1003	589	9.41	907	6.62	665
	RL	946	433	6.92	613	4.48	731

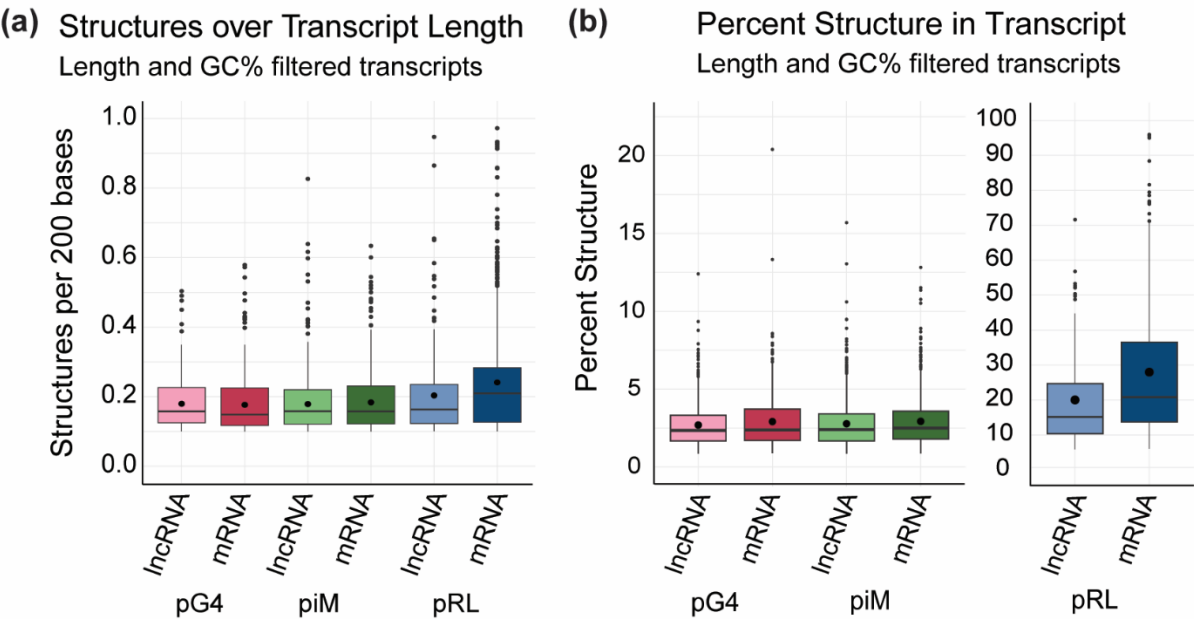


Figure S2. Unfiltered structure density and percent of transcript data for similar length and GC% lncRNA and mRNA dataset. (a) Box plot showing the distribution of the number of structures over transcript length for each type of predicted structure in lncRNA and mRNA with outliers. (b) Box plot showing the distribution of the total percentage of nucleotides in the transcript involved in structure formation for each structure type in lncRNA and mRNA including outliers.

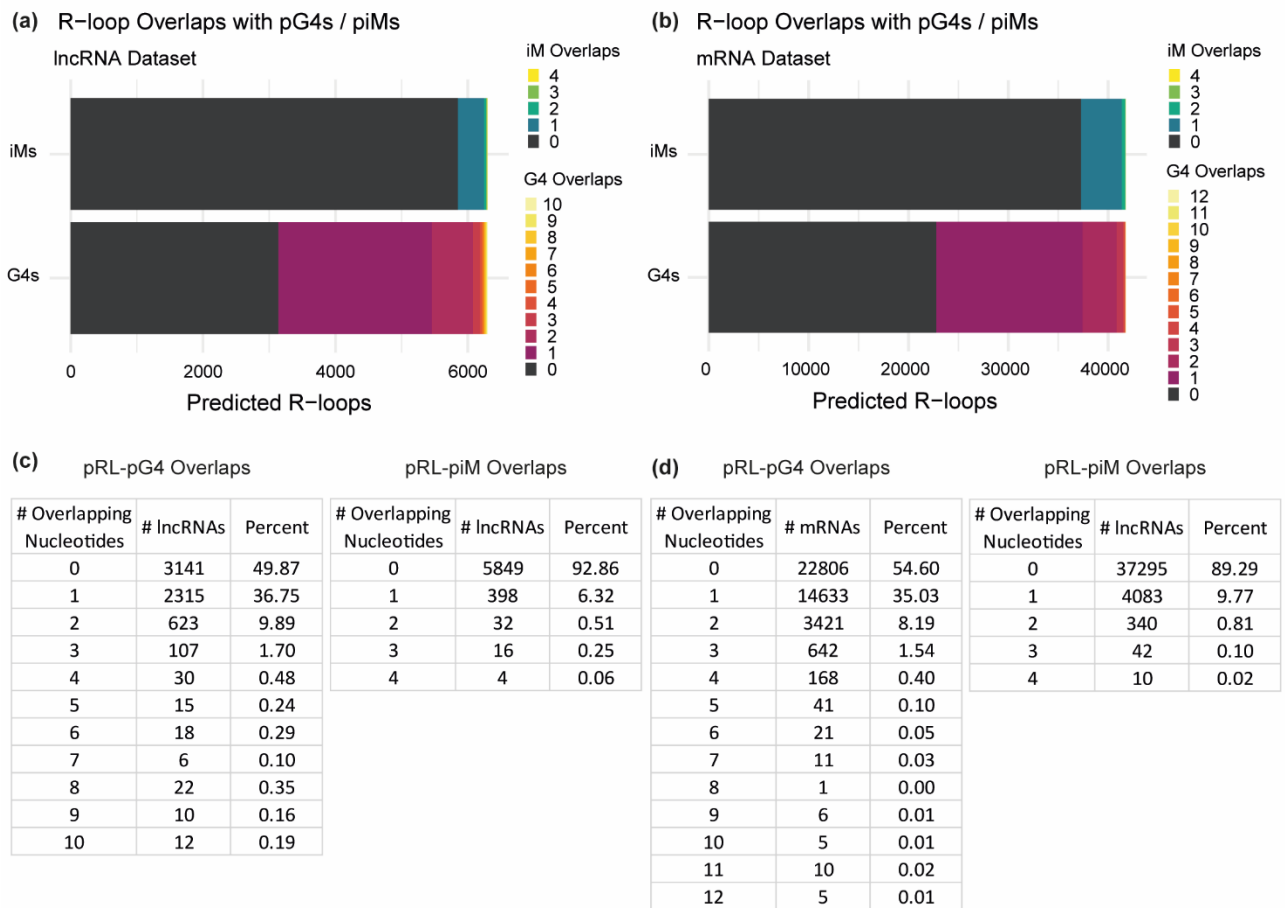


Figure S3. Sequence overlaps between predicted R-loops and predicted G4s and iMs. Stacked bar graphs showing pRL sequence overlaps with piM and pG4 sequences in **(a)** lncRNA and **(b)** mRNA. Colour legends represent the number of overlapping nucleotides from grey (lowest) to yellow (highest). Table with number and percentage of RNAs with given number of overlapping nucleotides between pRL and pG4/piM structures for **(c)** lncRNAs and **(d)** mRNAs.

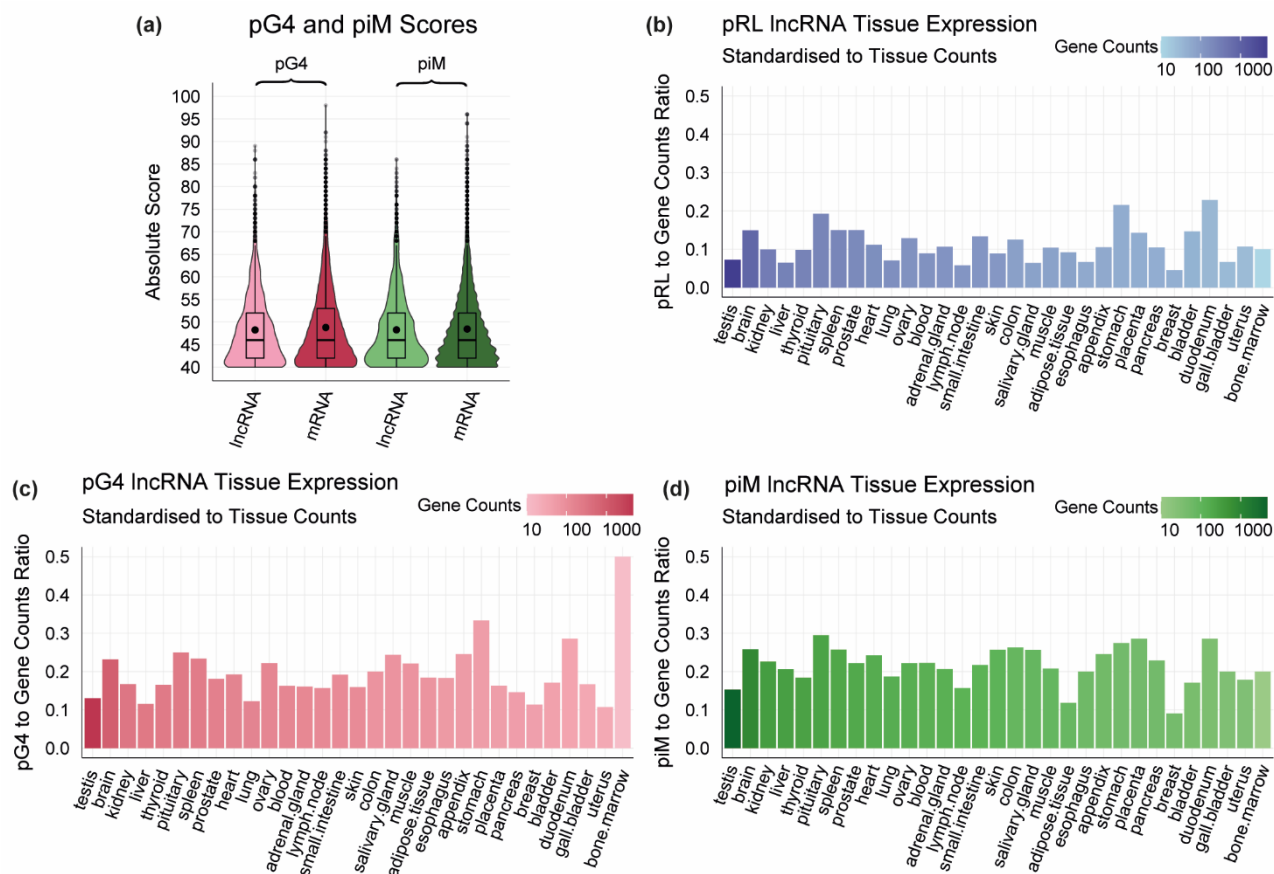


Figure S4. Score and tissue expression data for pG4s and piMs in lncRNA. (a) Violin + boxplot showing distribution of $|\text{score}|$ values for pG4 and piM structures in lncRNA and mRNA. Ratio of number of lncRNA genes with (b) pRL, (c) pG4, or (d) piM structures in a given tissue standardized to the number of genes in the tissue from the initial tissue dataset. Colour gradient indicates the number of genes in the tissue, light = low, dark = high. Tissues with less than 10 genes in the initial dataset were removed.

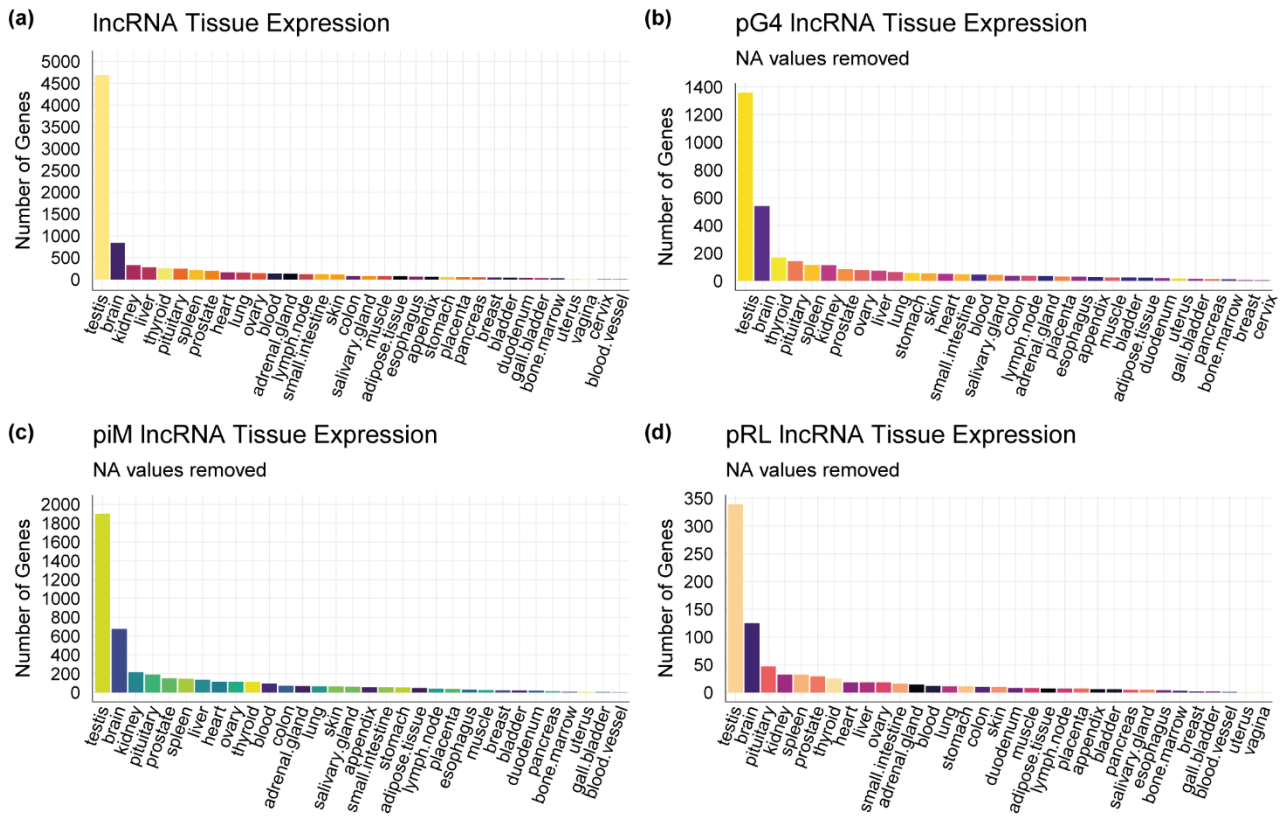


Figure S5. IncRNA tissue expression. (a) Total number of IncRNA genes per tissue in the input dataset without predicted structures. Total number of IncRNAs genes per tissue with predicted (b) G4, (c) iM, (d) and RL structures.

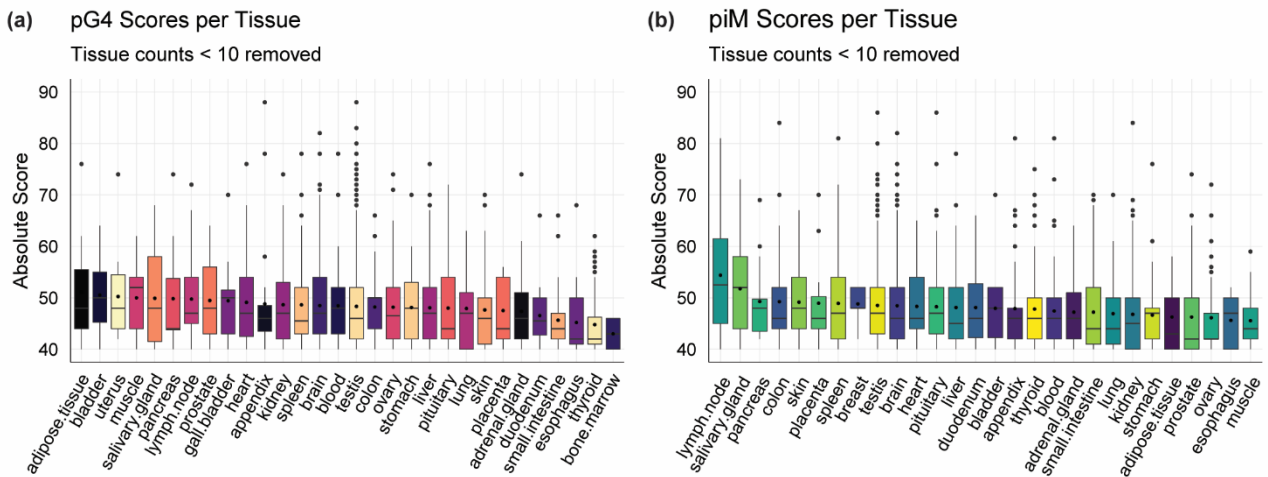


Figure S6. Score distribution in different tissues. Boxplots showing score distributions in each IncRNA tissue for (a) pG4 structures and (b) piM structures. Tissues with less than 10 genes in the initial dataset were removed.

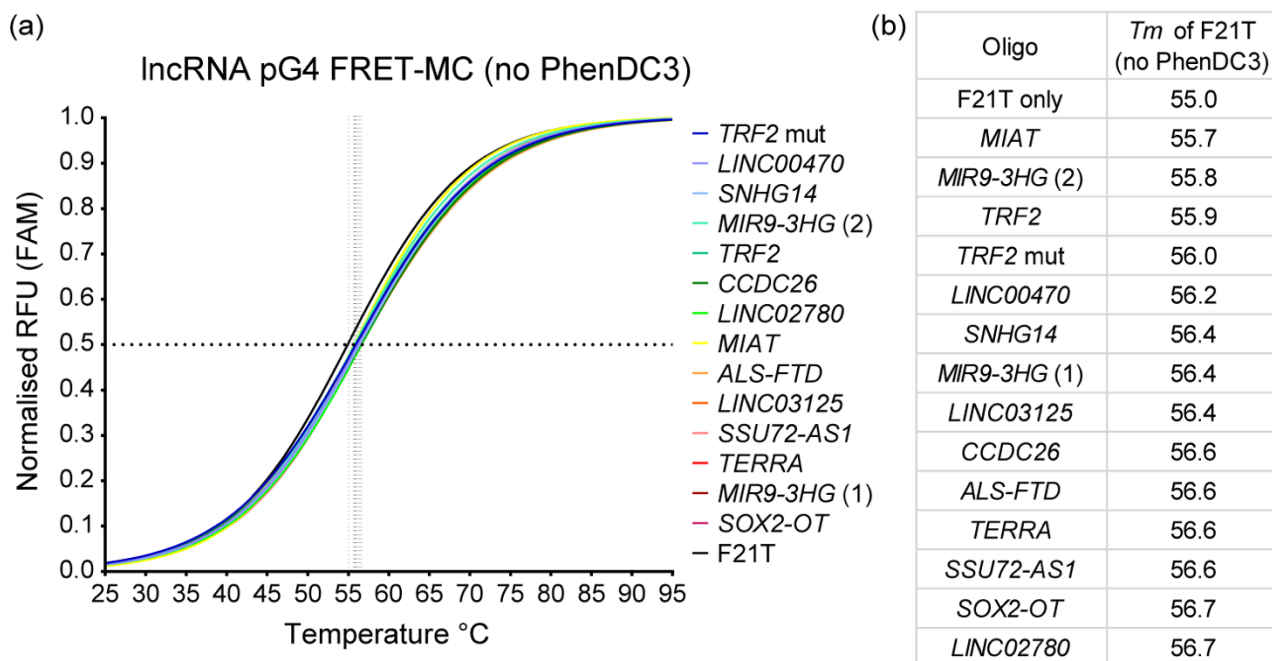


Figure S7. FRET-MC control data. (a) Non-linear fit of normalized F21T melt curves \pm test RNA pG4 competitor oligos without PhenDC3. Horizontal dotted line at $y = 0.5$, vertical dotted lines at T_m of F21T in each sample. (b) Table of T_m of F21T in all samples from FRET-MC without PhenDC3.

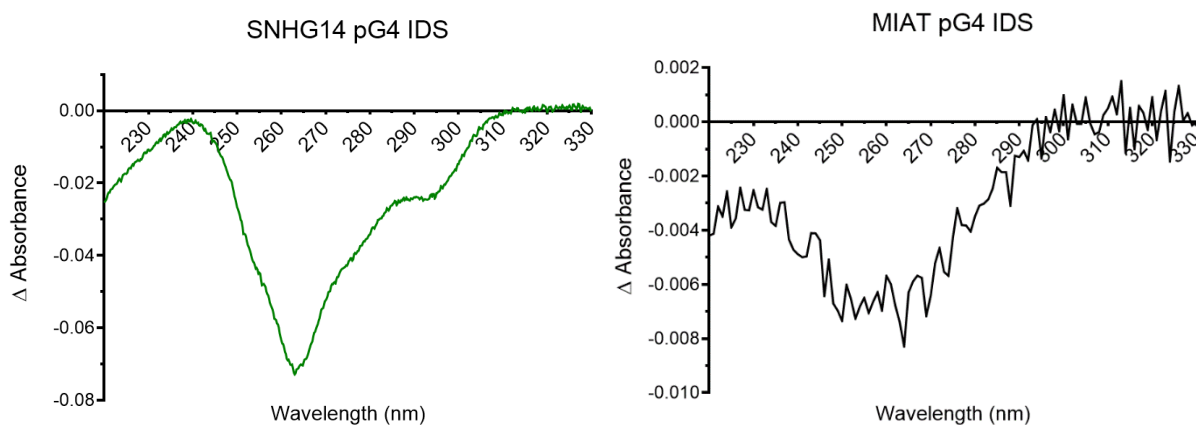


Figure S8. Raw IDS data for SNHG14 and MIAT pG4 lncRNA sequences at 37 °C.

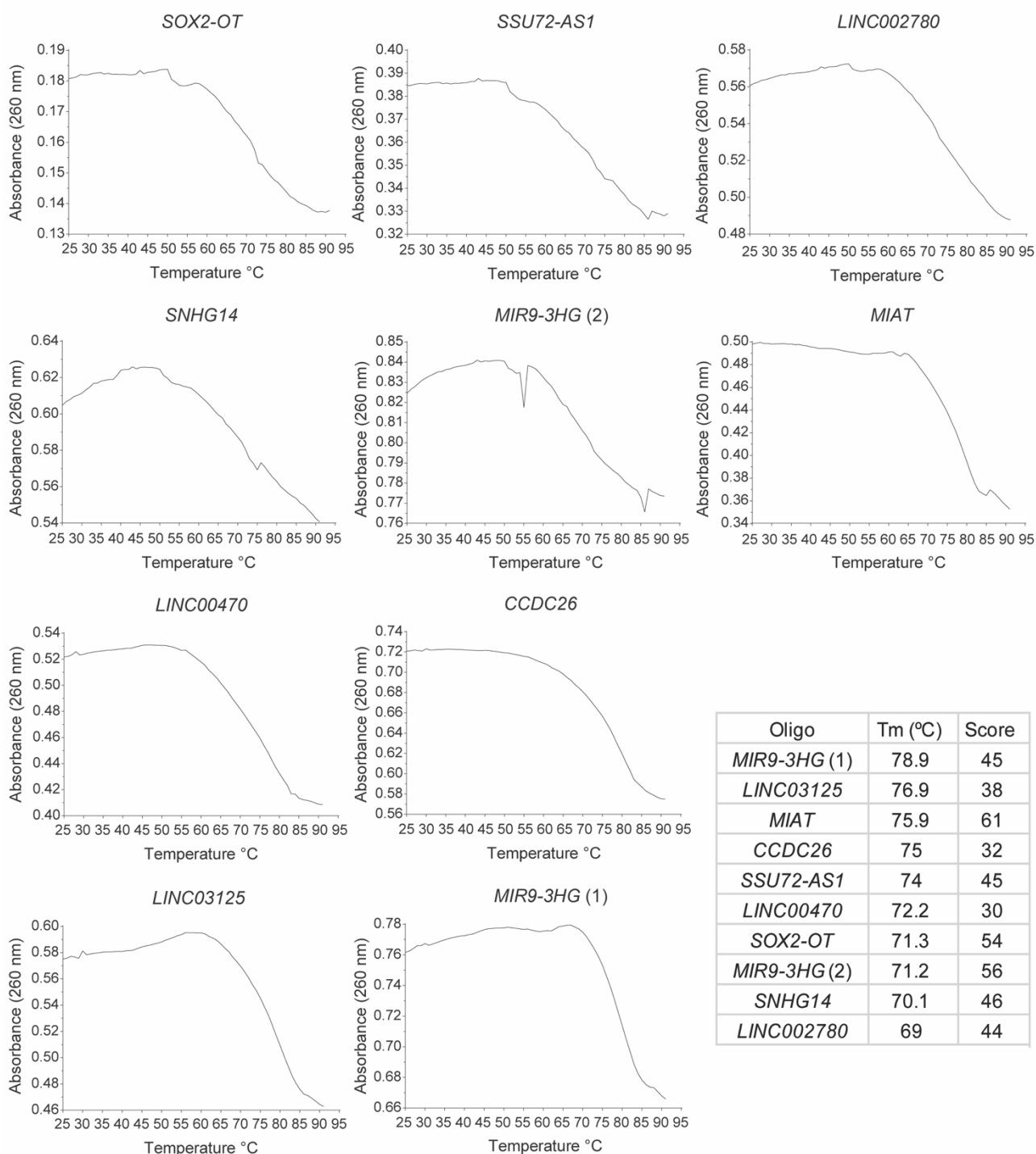


Figure S9. UV melt curve raw data and table showing T_m for each pG4 lncRNA oligo tested.

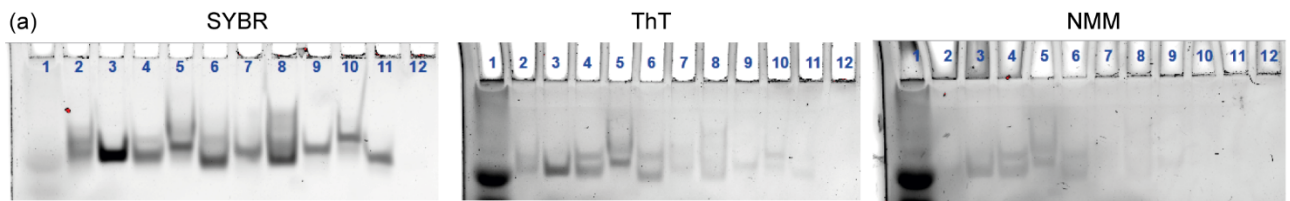


Figure S10. Polyacrylamide gel for pG4 lncRNA. 20% Native PAGE for lncRNA pG4 samples stained with (a) SYBR Safe, (b) ThT, or (c) NMM. Lanes – 1: Loading dye, 2: *SNHG14*, 3: *LINC02780*, 4: *SSU72-AS1*, 5: *SOX2-OT*, 6: *MIR9-3HG* (2), 7: *MIAT*, 8: *LINC00470*, 9: *CCDC26*, 10: *LINC03125*, 11: *MIR9-3HG* (1).

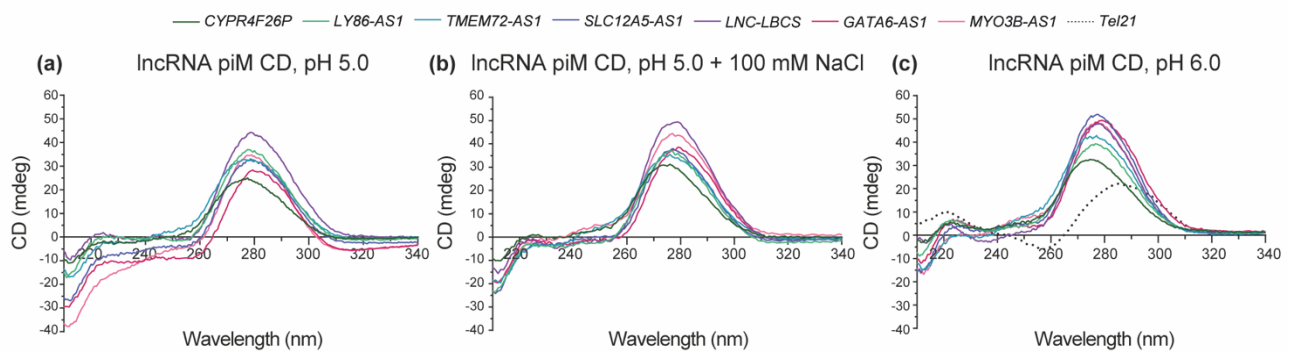


Figure S11. CD of lncRNA piM (Z-RNA) sequences. (a) 10 mM NaCaco, pH 5.0, (b) 10 mM NaCaco + 100 mM NaCl, pH 5.0, (c) 10 mM NaCaco, pH 6.0.

Table S2. pG4 sequences selected for biophysical analysis that are present in multiple lncRNAs. Name by which the sequence is referred to in the main manuscript is in bold.

5' - GGGAGGCUGAGGUGGG -3'					
LINC02780	<i>AC021148.3</i>	<i>AC130895.1</i>	<i>AZIN1-AS1</i>	<i>LINC00937</i>	<i>SHANK2-AS3</i>
<i>A1BG-AS1</i>	<i>AC022558.3</i>	<i>AC135012.1</i>	<i>BMP8B-AS1</i>	<i>LINC01476</i>	<i>SNHG5</i>
<i>AC003681.1</i>	<i>AC022596.2</i>	<i>AC138904.3</i>	<i>BSN-DT</i>	<i>LINC01511</i>	<i>SPACA6P-AS</i>
<i>AC004039.1</i>	<i>AC023347.1</i>	<i>AL022323.2</i>	<i>C12orf77</i>	<i>LINC01695</i>	<i>THAP9-AS1</i>
<i>AC004494.1</i>	<i>AC024588.1</i>	<i>AL022341.2</i>	<i>C21orf62-AS1</i>	<i>LINC02018</i>	<i>TMED2-DT</i>
<i>AC005005.4</i>	<i>AC026362.2</i>	<i>AL023882.1</i>	<i>CCDC15-DT</i>	<i>LINC02021</i>	<i>TMEM161B-AS1</i>
<i>AC006115.2</i>	<i>AC027644.3</i>	<i>AL109615.4</i>	<i>CTBP1-AS</i>	<i>LINC02595</i>	<i>TMEM254-AS1</i>
<i>AC006504.5</i>	<i>AC041005.1</i>	<i>AL157931.1</i>	<i>DENND6A-DT</i>	<i>LINC02642</i>	<i>TMEM51-AS1</i>
<i>AC007614.1</i>	<i>AC067863.1</i>	<i>AL359220.1</i>	<i>DLGAP1-AS1</i>	<i>MCPH1-AS1</i>	<i>TSIX</i>

<i>AC007922.4</i>	<i>AC067930.2</i>	<i>AL390071.1</i>	<i>EFCAB6-AS1</i>	<i>MIR3976HG</i>	<i>WAC-AS1</i>
<i>AC008543.5</i>	<i>AC073648.7</i>	<i>AL606753.2</i>	<i>ENTPD1-AS1</i>	<i>MKLN1-AS</i>	<i>XACT</i>
<i>AC008738.4</i>	<i>AC073957.3</i>	<i>AL627171.4</i>	<i>FAM153CP</i>	<i>MMP25-AS1</i>	<i>XIST</i>
<i>AC009318.1</i>	<i>AC097376.3</i>	<i>AL731559.1</i>	<i>HCG11</i>	<i>NCBP2-AS1</i>	<i>Z84486.1</i>
<i>AC010618.3</i>	<i>AC104024.2</i>	<i>AL732509.1</i>	<i>HELLPAR</i>	<i>PART1</i>	<i>Z95114.1</i>
<i>AC010733.1</i>	<i>AC107398.3</i>	<i>AP001020.1</i>	<i>LINC00461</i>	<i>PINK1-AS</i>	<i>Z98885.3</i>
<i>AC010907.2</i>	<i>AC109635.6</i>	<i>AP002518.1</i>	<i>LINC00467</i>	<i>PRR34</i>	
<i>AC016687.2</i>	<i>AC116158.2</i>	<i>AP002847.1</i>	<i>LINC00685</i>	<i>PSMA3-AS1</i>	
<i>AC021078.1</i>	<i>AC116914.2</i>	<i>AP003721.1</i>	<i>LINC00836</i>	<i>RFPLIS</i>	
5' -GGGCGUGGUGGCGGG-3'					
<i>SSU72-AS1</i>	<i>AC025580.3</i>	<i>AL009176.1</i>	<i>AL354892.3</i>	<i>CPB2-AS1</i>	<i>MKLN1-AS</i>
<i>AC004528.1</i>	<i>AC034229.1</i>	<i>AL022311.1</i>	<i>AL355312.2</i>	<i>DDR1-DT</i>	<i>SLC16A1-AS1</i>
<i>AC004941.1</i>	<i>AC073320.2</i>	<i>AL022322.2</i>	<i>AL356495.1</i>	<i>FAM153CP</i>	<i>SNHG14</i>
<i>AC005616.1</i>	<i>AC080080.2</i>	<i>AL022323.2</i>	<i>AL356756.1</i>	<i>GEMIN7-AS1</i>	<i>UCKL1-AS1</i>
<i>AC007608.3</i>	<i>AC090559.1</i>	<i>AL022344.1</i>	<i>AL390860.1</i>	<i>HELLPAR</i>	<i>XACT</i>
<i>AC007613.1</i>	<i>AC092447.5</i>	<i>AL031595.3</i>	<i>AL590666.3</i>	<i>KDM5C-IT1</i>	<i>Z95114.1</i>
<i>AC009403.2</i>	<i>AC093028.1</i>	<i>AL078459.1</i>	<i>AL591163.1</i>	<i>LIN28B-AS1</i>	<i>Z95114.2</i>
<i>AC016705.2</i>	<i>AC109322.2</i>	<i>AL133410.1</i>	<i>AL592211.1</i>	<i>LINC00887</i>	
<i>AC016747.1</i>	<i>AC112504.1</i>	<i>AL137783.1</i>	<i>AL928654.1</i>	<i>LINC00963</i>	
<i>AC020659.1</i>	<i>AC132938.3</i>	<i>AL139317.5</i>	<i>AP001977.1</i>	<i>LINC01285</i>	
<i>AC020916.1</i>	<i>AC132938.7</i>	<i>AL157886.1</i>	<i>CCDC26</i>	<i>LINC02840</i>	
<i>AC020928.2</i>	<i>AL008628.1</i>	<i>AL353708.1</i>	<i>CELF2-DT</i>	<i>LSINCT5</i>	
<i>AL645728.1</i>	<i>AC025580.3</i>	<i>AL009176.1</i>	<i>AL354892.3</i>	<i>CPB2-AS1</i>	

Table S3. Occurrence of the top 10 most prevalent lncRNA pG4 sequences in comparison to all pG4 motifs identified in lncRNA. Score = G4-iM Grinder Score. SINE = short interspersed nuclear element, LINE = long interspersed nuclear element, LTR = long terminal repeat.

Top 10 most prevalent lncRNA G4 sequences					
Sequence	# of occurrences	% in repeats	Repeat Class(es)	# unique transcripts	# unique genes
GGGAGGCUGAGGCAGGAG	560	54.1	SINE, LINE, LTR	508	358
GGAGGCUGAGGCAGGAG	365	57.6	SINE, LINE, LTR	344	233
GGGAGGCUGAGGUGGG	145	29.7	SINE	140	105
GGGCGUGGUGGCGGG	91	46.3	SINE	90	67
GGGAGGCCGAGGCGGG	81	33.8	SINE	80	65
GGGAGGCUGAGGCGGG	67	34.4	SINE	67	61
GGAGGCUGAGGUGGG	66	34.6	SINE	65	52
GGGAGGCCGAGGUGGG	86	20.0	SINE	86	45
GAGGCUGAGGCAGGAG	69	55.0	SINE, LTR	68	40
GGGAGGCUGAGGCAGGAGAAU GGCGUGAACCCGGGAGGCG	48	36.7	SINE	48	30
All lncRNA pG4s (Score ≥ 20 , MinNRuns ≥ 3)	82,284	20.1	-	28,313	11,262
All lncRNA pG4s (Score ≥ 40 , MinNRuns ≥ 4)	6,645	34.2	-	5,023	2,964
All mRNA pG4s (Score ≥ 40 , MinNRuns ≥ 4)	29,049	25.8	-	20,046	7,912