

1 **Single position substitution of hairpin pyrrole-imidazole polyamides**
2 **imparts distinct DNA-binding profiles across the human genome**

3 Paul B. Finn^{1¶, #a}, Devesh Bhimsaria^{2¶}, Asfa Ali³, Asuka Eguchi⁴, Aseem Z. Ansari⁵, Peter B. Dervan^{1*},

4 ¹ Division of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena,
5 California, United States of America

6 ² Bio Informatics, Jaipur, Rajasthan, India

7 ³ Department of Molecular Genetics, University of Texas Southwestern Medical Center, Dallas,
8 Texas, United States of America

9 ⁴ Department of Microbiology and Immunology, Stanford University, Stanford, California, United
10 States of America

11 ⁵ Department of Chemical Biology & Therapeutics, St. Jude Children's Research Hospital, Memphis,
12 Tennessee, United States of America

13 ^{#a} Current Address: Department of Bioengineering, Stanford University, Stanford, California, United
14 States of America

15
16 * Corresponding author

17
18 E-mail: dervan@caltech.edu (PBD)

19
20
21 [¶] These authors contributed equally to this work.

22 **ABSTRACT**

23 Regulating desired loci in the genome with sequence-specific DNA-binding molecules is a
24 major goal for the development of precision medicine. Pyrrole–imidazole (Py–Im) polyamides
25 are synthetic molecules that can be rationally designed to target specific DNA sequences to
26 both disrupt and recruit transcriptional machinery. While *in vitro* binding has been extensively
27 studied, *in vivo* effects are often difficult to predict using current models of DNA binding.
28 Determining the impact of genomic architecture and the local chromatin landscape on
29 polyamide-DNA sequence specificity remains an unresolved question that impedes their
30 effective deployment *in vivo*. In this report we identified polyamide–DNA interaction sites
31 across the entire genome, by covalently crosslinking and capturing these events in the nuclei
32 of human LNCaP cells. This method, termed COSMIC-seq, confirms the ability of hairpin-
33 polyamides, with similar architectures but differing at a single ring position, to retain *in vitro*
34 specificities and display distinct genome-wide binding profiles. These results underpin the
35 development of Py-Im polyamides as DNA-targeting molecules that mediate their regulatory
36 or remedial functions at desired genomic loci.

37 INTRODUCTION

38 Regulating genomic architecture and activity with sequence-specific synthetic DNA
39 binding molecules is a long-standing goal at the interface of chemistry, biology and medicine.
40 Small molecules that selectively target desired genomic loci could be harnessed to regulate
41 critical gene networks. The greatest success in designing small molecules with programmable
42 DNA-binding specificity has been with pyrrole-imidazole (Py-Im) polyamides [1–8]. Pyrrole-
43 imidazole (Py-Im) polyamides are synthetic DNA-binding oligomers with high sequence
44 specificity and affinity [7]. An oligomer, comprising a modular set of aromatic pyrrole and
45 imidazole amino acids linked in series by a central aliphatic γ -aminobutyric acid (GABA) ‘turn’
46 unit, fold into a hairpin structure in the minor groove of DNA and afford binding affinities and
47 specificities comparable to natural transcription factors [3,7]. Sequence specificity is
48 programmed through side-by-side pairs of the Py and Im subunits that “read” the steric and
49 hydrogen bonding patterns presented by the edges of the four Watson-Crick base pairs on
50 the floor of the minor groove [5]. DNase I footprinting titrations and other *in vitro* methods
51 have extensively characterized the binding affinity and specificity of these molecules [3,6,7,9].
52 An Im/Py pair binds G•C; Py/Im binds C•G, and Py/Py pairs both bind A•T and T•A (denoted
53 as W) [1,2]. Py-Im polyamide binding in the minor groove induces allosteric changes to DNA,
54 widening the minor groove and narrowing the major groove [10–12]. Polyamide-DNA binding
55 is sufficient to disrupt protein-DNA interfaces, including DNA interactions made by
56 transcription factors and the transcriptional machinery [13–15]. Additionally, polyamides can
57 function as sequence-specific synthetic cofactors through allosteric DNA modulation to
58 enhance the assembly of protein-DNA complexes [12]. Py-Im polyamides are cell permeable,
59 localize to the nucleus in live cells and are non-genotoxic [16–18] failing to activate canonical
60 DNA damage response or significantly alter cell cycle distribution [19].

61 The identification of new mechanistic insights into Py-Im polyamide activity have underlined
62 the importance of mapping polyamide binding to chromatin [15,18,19]. Polyamide binding in
63 the more complex cellular environment presents a formidable challenge since chromatin DNA
64 has varying degrees of accessibility. Sequence specific access by Py-Im polyamides to the
65 nucleosome core particle (NCP) has been demonstrated *in vitro* and with x-ray crystal
66 structures of NCP•polyamide complexes [20–22]. However, the extent to which chromatin
67 states influence polyamide binding to its cognate sites remains a long-standing question. The
68 lack of clarity on the parameters that govern genome-wide binding of polyamides greatly
69 impedes the deployment of this class of molecules to regulate cell fate-defining and disease-
70 causing gene networks *in vivo*.

71 We report here the genome-wide binding profiles of two Py-Im polyamides **1** and **2**, of
72 identical architecture (8-ring hairpin) that differ at a single aromatic ring position in cellular
73 nuclei using COSMIC-seq ('crosslinking of small molecules for isolation of chromatin with next-
74 generation sequencing), **Fig 1** [23,24]. COSMIC-seq employs a tripartite conjugate composed
75 of the DNA-binding ligand attached to a biotin affinity handle and a psoralen photocrosslinker.
76 Genome-wide binding of these tripartite molecules is captured by photo-induced crosslinking
77 followed by biotin-enabled enrichment and unbiased NGS sequencing of the conjugated
78 genomic loci [23,24]. The ability to induce rapid crosslinking at the desired time point
79 distinguishes COSMIC-seq from continuous and uncontrolled alkylation-dependent DNA
80 conjugations that have been used to query genome-wide binding of polyamides [25].
81 COSMIC-seq also differs from Chem-seq approaches that use ligands for protein complexes
82 that are associated with the genome [26]. Previously, COSMIC-seq was utilized to access
83 genome-wide binding of two structurally distinct Py-Im polyamides (hairpin vs linear) that code
84 for very different sequences [24]. An 8-ring hairpin Py-Im polyamide (TpPyPyIm- γ -PyImPyPy-
85 β -Dp) binds 6 bp of DNA (5'-WTWCGW-3') [27], whereas a linear polyamide (ImPy- β -ImPy-
86 β -Im- β -Dp) binds 9 bp of purine rich DNA (5'-AAGAAGAAG-3') [28–31]. While such a dramatic
87 difference in target sequence composition leads to distinct genome-wide binding profiles, we
88 wondered how a more challenging single position change (CH to N:) within one ring of an 8-
89 ring hairpin would affect genomic occupancy. In this study we applied COSMIC-seq to
90 determine if two polyamides of identical size and architecture, hairpins **1** and **2** which code for
91 6 base pair sites differing by one base pair position 5'-WGWWCW-3' and 5'-WGGWCW-3',
92 respectively, can display distinct genomic binding occupancy on chromatin. These
93 experiments provide a more stringent test of genome-wide binding properties of hairpin
94 polyamides in a chromatin environment for application as precision-targeting molecules.

95

96 **Fig 1. Trifunctional Py-Im polyamide conjugates 1 and 2.** (A) Chemical structure of hairpin
97 Py-Im polyamides **1** and **2** which differ by one atom, shown in red, and (B) the corresponding
98 predicted target sequences based on the pairing rules. Py-Im polyamide **1** targets the DNA
99 sequence 5'-WGWWCW-3' and Py-Im polyamide **2** targets 5'-WGGWCW-3'. Open and filled
100 circles represent N-methylpyrrole (Py) and N-methylimidazole (Im), respectively. The N-
101 acetylated (R)- γ -aminobutyric acid turn residue is shown as a semicircle, and psoralen and
102 biotin are denoted by P and B, respectively.

103

104 MATERIALS AND METHODS

105 **Materials**

106 Chemicals and solvents were purchased from standard chemical suppliers and used without
107 further purification. (R)-2,4-Fmoc-Dab(Boc)-OH (α -amino-GABA turn) was purchased from
108 Peptides International. Monomers were synthesized as previously described [32]. Kaiser
109 oxime resin (100-200 mesh) and benzotriazole-1-yl-oxy-trispyrrolidinophosphonium
110 hexafluorophosphate (PyBOP) were purchased from Novabiochem. 2-Chlorotriyl chloride
111 resin was purchased from Aapptec. Preparative HPLC purification was performed on an
112 Agilent 1200 Series instrument equipped with a Phenomenex Gemini preparative column (250
113 x 21.2 mm, 5 μ m) with the mobile phase consisting of a gradient of acetonitrile (CH₃CN) in 0.1%
114 aqueous trifluoroacetic acid (TFA). Polyamide concentrations were measured by UV/Vis
115 spectroscopy in distilled and deionized water (ddH₂O) with a molar extinction coefficient of
116 8650 M⁻¹ cm⁻¹ at 310 nm for each *N*-methylpyrrole (Py) and *N*-methylimidazole (Im) and 11,800
117 M⁻¹ cm⁻¹ for the psoralen/biotin derivative **3** [33,34]. Analytical HPLC analysis was conducted
118 on a Beckman Gold instrument equipped with a Phenomenex Gemini analytical column (250
119 x 4.6 mm, 5 μ m), a diode array detector, and the mobile phase consisting of a gradient of
120 acetonitrile in 0.1% aqueous TFA. Matrix-assisted, LASER desorption/ionization time-of-flight
121 (MALDI-TOF) mass spectrometry was performed on an Autoflex MALDI TOF/TOF (Bruker)
122 using α -cyano-4-hydroxycinnamic acid matrix. Oligonucleotides were purchased from
123 Integrated DNA Technologies Inc. All sequencing samples were processed as single read
124 (50 bp) sequencing runs at the California Institute of Technology Millard and Muriel Jacobs
125 Genetics and Genomics Laboratory on an Illumina HiSeq 2500 Genome Analyzer.

126 **Chemical synthesis**

127 Polyamides **1A** and **2A** were synthesized on solid support (Kaiser oxime resin, 100-200 mesh),
128 using microwave-assisted PyBOP coupling conditions with *N*-methylpyrrole (Py), *N*-
129 methylimidazole (Im) amino acid monomers and dimers (**5a** & **5b**) as previously described,
130 **S1A Fig** [34]. Polyamides were cleaved from resin with neat 3,3'-diamino-*N*-
131 methylpropylamine (60 °C, 5 min, μ W), precipitated with diethyl ether at -20 °C, re-dissolved
132 in 20 - 30% (v/v) CH₃CN/H₂O (0.1% TFA), and purified by reverse-phase preparative HPLC.
133 Fractions that showed clean polyamide without contaminants were frozen in liquid nitrogen
134 and lyophilized to dryness as a white-yellow solid. The identity and purity were confirmed by
135 MALDI-TOF mass spectrometry and analytical HPLC. The observed mass for **1A**
136 (C₅₉H₇₅N₂₂O₁₀) is 1251.78 (calculated 1251.60) and for **2A** (C₅₈H₇₄N₂₃O₁₀) is 1252.75
137 (calculated 1252.60).

138 The psoralen-biotin peptide **3** was synthesized by manual Fmoc solid-phase synthesis on 2-
139 chlorotrityl chloride resin by standard procedures, **S1B Fig** [23]. Coupling and deprotection
140 were performed at room temperature for 1 h and 15 min, respectively. Briefly, Fmoc-protected
141 amino acids or polyethylene glycol (PEG) linkers were activated with HATU and HOAt in the
142 presence of *N,N*-diisopropylethylamine (DIPEA) in dimethylformamide (DMF) (or DMSO/DMF)
143 and deprotection of the Fmoc group was achieved with 20% piperidine in DMF. Cleavage
144 from resin was achieved with a solution of 95% (v/v) TFA, 2.5% (v/v) H₂O, and 2.5% (v/v)
145 triisopropylsilane and purified by reverse-phase preparative HPLC, lyophilized to dryness as
146 a white powder and protected from light. The identity and purity were confirmed by MALDI-
147 TOF mass spectrometry and analytical HPLC. The observed mass for **3** (C₄₃H₆₁N₆O₁₅S) is
148 933.36 (calculated 933.39).

149 Polyamide-peptide conjugates **1** and **2** were synthesized by solution phase peptide coupling
150 conditions and protected from light, **S1C Fig**. Peptide acid **3** (1 equiv.) was pre-activated for
151 5 min at room temperature with a solution of HATU/HOAt/DIPEA (3:3:6 equiv.) in DMF.
152 Polyamide **1A** or **2A** was added (1 - 1.5 equiv.), and the coupling was allowed to proceed for
153 30 - 60 minutes until all of **3** was consumed as determined by analytical HPLC. The
154 polyamide-peptide conjugates were purified by reverse-phase HPLC and lyophilized to
155 dryness. The identity and purity were confirmed by MALDI-TOF mass spectrometry and
156 analytical HPLC. The observed mass for **1** (C₁₀₂H₁₃₃N₂₈O₂₄S) is 2166.26 (calculated 2165.98)
157 and for **2** (C₁₀₁H₁₃₂N₂₉O₂₄S) is 2167.23 (calculated 2166.97).

158 **Cognate site identification**

159 High-throughput cognate binding sites were identified for the polyamides **1** and **2** using SELEX
160 method [35]. A DNA library with a central randomized 20-bp region and flanked by constant
161 sequences (~10¹² possible sequences, Integrated DNA Technologies) was used for PCR
162 amplification. Polyamide conjugates **1** and **2** at a range of concentrations (5 nM and 50 nM)
163 were added to 100 nM of DNA library in binding buffer [1× PBS (pH 7.6), 50 ng/ μL poly(dI-
164 dC)] and incubated for 1 h at room temperature. Enrichment of the compound-DNA complexes
165 was performed using streptavidin-coated magnetic beads (Dynabeads, Invitrogen) following
166 manufacturer's protocol. To remove unbound DNA, three washes were done after the capture,
167 with 100 μL ice-cold binding buffer. Beads were resuspended in PCR master mix (EconoTaq
168 PLUS 2× Master Mix, Lucigen), the DNA was amplified for 15 cycles and purified (QIAGEN).
169 Three rounds of selection were performed (DNA was quantified by absorbance at 260 nm
170 before each round of binding). An additional round of PCR was performed after completion of
171 three rounds of selection, to incorporate Illumina sequencing adapters and a unique 6-bp
172 barcode for multiplexing. The starting library was also barcoded and sequenced. Samples

173 were sequenced on an Illumina HiSeq 2500 at the Millard and Muriel Jacobs Genetics and
174 Genomics Laboratory in California Institute of Technology.

175 **Cell culture conditions**

176 LNCaP cells were maintained in RPMI 1640 (Invitrogen) with 10% Fetal Bovine Serum (FBS,
177 Irvine Scientific) at 37°C under 5% CO₂. LNCaP cells were purchased from ATCC (Manassas,
178 VA, USA).

179 **Crosslinking of small molecules for isolation of chromatin** 180 **with next-generation sequencing**

181 COSMIC-seq was performed in LNCaP nuclei, as previously described [23]. LNCaP cells
182 (~2.5 x 10⁷) were washed twice with cold PBS then resuspended in cold lysis buffer (RSB +
183 0.1% IGEPAL CA-630, 2.5 x 10⁷ cells/250 µL), incubated on ice for 5 min then centrifuged
184 immediately at 130 x g for 10 min at 4 °C. Nuclei were resuspended in binding buffer [10 mM
185 Tris HCl (pH 8.0), 5 mM MgCl₂, 1 mM DTT, 0.3 M KCl, 0.1 M PMSF, 0.1 M benzamidine, 0.1
186 M pepstatin A, 10% glycerol] and treated with psoralen-biotin conjugated polyamide **1** or **2**
187 (0.4 µM and 4 µM, 0.1% DMSO final concentration) for 1 h at 4 °C in the dark. Nuclei were
188 irradiated for 30 min with a UV lamp (2.4 µW/cm²; CalSun) through a Pyrex filter, centrifuged
189 at 500 x g and re-suspended in COSMIC buffer [20 mM Tris·Cl (pH 8.1), 2 mM EDTA, 150
190 mM NaCl, 1mM PMSF, 1mM benzamidine, 1.5 µM pepstatin, 1% Triton X-100, 0.1% SDS].
191 Samples were sonicated at 3 °C for 36 min with a cycle of 10 s ON and 10 s OFF, at HIGH
192 setting (Bioruptor Plus, Diagenode). Samples were centrifuged 10 min at 12,000 x g and 10%
193 of the sample was saved as input DNA and stored at -80 °C until reversal of cross-linking.
194 The rest of the sample was used for the affinity purification (AP). Streptavidin-coated magnetic
195 beads (100 µL per sample, Dynabeads MyOne C1) were washed in COSMIC buffer and
196 incubated with AP samples for 16 h at 4°C. All washes were performed at room temperature
197 unless otherwise noted. For **1** and **2**, **1A** and **2A** were added (5 µM), respectively, in the
198 washes. Samples were washed twice with COSMIC buffer (once 12 h and once 4 h). Samples
199 were then washed once with washing buffer 1 [10 mM Tris·Cl (pH 8.0), 1 mM EDTA, 3% (v/v)
200 SDS], once with washing buffer 2 [10 mM Tris·Cl (pH 8.0), 250 mM LiCl, 1 mM EDTA, 0.5%
201 Nonidet P-40, 1% sodium deoxycholate], twice with freshly prepared washing buffer 3 [4 M
202 urea, 10 mM Tris·Cl (pH 7.5), 1 mM EDTA, 0.1% Nonidet P-40], and twice with TE buffer [10
203 mM Tris·Cl (pH 8.0), 1 mM EDTA]. Samples were re-suspended in TE and labelled as AP
204 DNA. Input and AP samples were re-suspended in cross-link reversal buffer [10 mM Tris (pH
205 7.6), 0.4 mM EDTA, 100 mM KOH]. Crosslinks were reversed, and DNA was eluted from

206 beads at the same time by heating samples for 30 min at 90 °C. Input and AP samples were
207 neutralized with 6N HCl, and incubated first with RNase A (0.2 µg/µL) for 1 h at 37 °C and
208 then with Proteinase K (0.2 µg/ µL) for 1 h at 55 °C. Samples were purified with the MinElute
209 PCR Purification Kit (Qiagen).

210 **CSI data analysis**

211 The reads from Illumina sequencing were de-multiplexed using the 6 bp barcodes and then
212 truncated to include only the 20 bp random portion of the library. On average, 1,031,000 reads
213 per barcode were obtained. The occurrence of every k-mer (8 mer) was counted using a
214 sliding window of size k. To correct for experimental biases and biases in the initial DNA
215 library, a standardized enrichment score was calculated by normalizing the counts of every
216 k-mer from the enriched CSI data (rounds 1, 2 or 3) to the expected number of counts in the
217 library with a fifth-order Markov model derived from the processed library (processed same
218 number of SELEX enrichment rounds without the polyamide as done for the polyamide)
219 [36,37]. The most enriched 8 bp sub-sequences were used to derive position weight matrix
220 (PWM) motifs using MEME [38,39]. Data files for mapped 20 bp reads and normalized 8 bp
221 sequences are available online (<https://ansarilab.biochem.wisc.edu/computation.html>).

222 **Sequence logos**

223 PWMs were derived from the 50 most enriched 8-mer sequences (ranked by enrichment) for
224 each polyamide, using MEME [38,39]. MEME was run with following parameters:

225 `-dna -mod anr -nmotifs 10 -minw 6 -maxw 8 -time 7200 -maxsize 60000 -revcomp`

226 **Specificity and energy landscapes**

227 Specificity and Energy Landscapes (SELs) display high-throughput protein-DNA binding data
228 (DNA–protein interactome or DPI) in the form of concentric rings [40–42]. The organization of
229 data in SEL is detailed in **S3 Fig**. SELs were generated from 8-mer enrichment files using the
230 target sequence for corresponding polyamide **1** (5'-WGWWCW-3') and **2** (5'-WGGWCW-3')
231 as seed motif. The software for generating SELs is made available online
232 (<https://ansarilab.biochem.wisc.edu/computation.html>).

233 **Genomescales: scoring *in vivo* bound sites with *in vitro*** 234 **data**

235 Genomescaples are generated by assigning *in vitro* CSI intensities (enrichment values) to
236 genomic regions. To generate CSI Genomescaples a sliding k-mer window was used to score
237 genomic regions and then plotted as a bar plot [40,42].

238 **Summation of sites model**

239 Summation of sites (SOS) model was used to predict DNA binding of polyamides in the human
240 genome, hg19 [23,24,43]: The SOS score was obtained by summing (or averaging) all k-mer
241 *in vitro* binding intensities (enrichment) obtained using a sliding k-mer window across a
242 genomic region. Data is displayed using genomic regions of 420 bp for SOS [23]. For SOS
243 predicted genomic loci, the whole human genome (hg19) was divided into 420 bp fragments
244 with the overlap of half (210 bp). These fragments were then sorted by the predicted binding
245 to polyamide **1** and **2** using SOS model. The top 1000 predicted peaks obtained were used as
246 final predicted peaks for further analysis.

247 **COSMIC-seq data analysis**

248 Sequencing reads were mapped to the human genome (hg19) with Bowtie (best -m 1) to yield
249 unique alignments. Bound regions/peaks were identified with SPP [24,43]. The data has
250 been deposited in the Gene Expression Omnibus (GEO) database, www.ncbi.nlm.nih.gov/geo
251 (accession no. GSE149367).

252 **Plotting tag density data, genomescape and SOS as** 253 **heatmaps**

254 To display data multiple heatmaps are shown using two types of genomic regions: the top
255 1000 COSMIC peaks, and the top 1000 genomic loci predicted to be the best binder using the
256 SOS model. These regions were scored using COSMIC-seq tag density for AP of a 10 Kbp
257 region surrounding the peak using HOMER `annotatePeaks.pl` command with arguments `-hist`
258 `as 25 bp`, SOS scores for 10 Kbp region surrounding the peak, and genomescales for 1 Kbp
259 [44]. Different coloring scales were used to display heatmaps by using a multiplication factor
260 of 10x for tag density and 100x for genomescales and SOS scoring.

261 **RESULTS**

262 **Polyamide Design**

263 COSMIC-seq was performed on two structurally identical hairpin polyamides, differing at a
264 single position (X = CH vs N:) on the second aromatic amino acid ring-pair, **Fig 1A**. A single
265 CH to N: position substitution changes the ring pair from a Py/Py to an Im/Py which invokes a
266 preference from an A•T or T•A to a G•C base pair, respectively, based on previously
267 determined pairing rules, **Fig 1B** [1,2,45]. Py-Im polyamide **1**, designed to target the
268 consensus androgen response element (ARE) half-site 5'-WGWWCW-3', has been shown to
269 regulate androgen receptor (AR) and glucocorticoid (GR) driven gene expression in cell
270 culture and suppress tumor growth *in vivo* [14,18,46]. Py-Im polyamide **2**, designed to target
271 the estrogen response element (ERE) consensus half site 5'-WGGWCW-3', was shown to
272 effect estrogen receptor-alpha (ER α)-driven gene expression *in vitro* and *in vivo* [47]. In this
273 study each, hairpin Py-Im polyamide is conjugated at the C-terminus with a psoralen and biotin
274 for enrichment connected via a linker (~36 Å extended) capable of sampling pyrimidine
275 proximal to the polyamide-binding site suitable for 2 + 2 photocycloaddition [23,24]. Because
276 the psoralen moiety crosslinks proximal pyrimidines (T in particular), we anticipate subtle bias
277 in the data, a contextual flanking sequence nuance adjacent to the core binding sites of each
278 polyamide. Py-Im polyamides were synthesized by Boc solid-phase synthesis, cleaved from
279 resin and conjugated to the psoralen-biotin moiety **3** (**S1 Fig**).

280

281 **Different sequence specificities conferred by a single** 282 **position substitution**

283 To comprehensively map *in vitro* binding characteristics of hairpin polyamide-conjugates **1**
284 and **2**, we performed solution-based Cognate Site Identifier (CSI) analysis, **Fig 2** [6,40,41,48].
285 Sequence specificity data was determined with next generation sequencing (NGS) by
286 solution-based enrichment methods (SELEX-seq) to assess polyamide-DNA binding [49].
287 These methods provide a comprehensive characterization of polyamide-DNA binding through
288 the sampling of a large sequence space (a dsDNA library bearing all $\sim 10^{12}$ sequence
289 permutations of a 20-bp site) using affinity purification coupled with massively parallel
290 sequencing [35,41]. This platform allows rapid, quantitative identification of the full spectrum
291 of polyamide binding sites of up to 20 bp in size, correlates well with solution-phase and
292 microarray platforms, and has been used to guide the refinement of general polyamide design
293 principles [9,24,31,40,41]. Py-Im polyamides **1** and **2** were incubated with a duplex
294 oligonucleotide library containing a randomized 20-mer region, and the bound and unbound
295 sequences were separated via affinity purification by streptavidin-coated magnetic beads, **Fig**
296 **2A**. Following each round of enrichment, sequences were PCR amplified, purified,
297 multiplexed, and subjected to massively parallel sequencing analysis. Computational analysis

298 was applied to enriched sequences to obtain binding site intensity values corresponding to all
299 8-mer DNA sequences, see methods [37].

300

301 **Fig 2. Detection of *in vitro* DNA binding of polyamides 1 and 2 via cognate site**
302 **identification (CSI) by SELEX-seq.** (A) Overview of CSI by SELEX-seq workflow. A
303 randomized 20 bp DNA library is incubated with biotinylated polyamide, DNA is enriched by
304 streptavidin-coated magnetic beads, PCR amplified and sequenced by NGS to obtain k-mers
305 (CSI enrichment) representing polyamide-DNA binding. Enrichment is displayed as a
306 histogram plot and high binding sequences are represented as a position weight matrix (PWM)
307 logo. Specificity and energy landscapes (SEs) are created to visualize the full spectrum of
308 DNA binding across all sequence permutations of an 8-mer binding site. (B) PWM logos for
309 polyamides **1** (*top*) and **2** (*bottom*). (C) Scatterplot comparison of *in vitro* DNA binding for **1** vs
310 **2**. CSI enrichment for 8-mers is plotted for sequences containing 5'-WGWWCW-3' (*red*) and
311 5'-WGGWCW-3' (*green*). (D) Comprehensive SEs for **1** (*top*) and **2** (*bottom*) using 5'-
312 WGWWCW-3' and 5'-WGGWCW-3' as seed motif, respectively, where W = A or T.

313 Polyamide-DNA binding motifs for each round of enrichment were identified by position weight
314 matrices (PWMs) using the top 50 enriched 8-mer sequences and displayed as sequence
315 logos, **Fig 2B** and **S2 Fig**, [38,39]. The highest information content for each polyamide is
316 found at a binding site width of six, verifying the binding site size expected when **1** and **2** are
317 bound in a fully ring-paired, hairpin configuration. The motifs generated are indicative of
318 polyamide-DNA binding consistent with the Py-Im pairing rules for both **1** and **2**, targeting the
319 sequences 5'-WGWWCW-3' and 5'-WGGWCW-3', respectively. A clear difference in the
320 sequence preference at the third position, corresponding to the CH to N: position substitution
321 of the second ring pair (Py/Py vs Im/Py), was detected. Additionally, both polyamides show
322 subtle differences in binding preference at the fourth and sixth positions revealing sensitivity
323 of binding energetics to changes in sequence context. Scatter plot comparison analysis of all
324 enriched 8-mer sequences indicates a preference for the consensus motif, polyamide **1**
325 prefers WGWWCW over WGGWCW, whereas polyamide **2** prefers WGGWCW over
326 WGWWCW, **Fig 2C** and **S3 Fig**. These results demonstrate that a single position (CH to N:)
327 modification of the aromatic amino acid ring of the polyamide core structure imparts a
328 significant change in the global *in vitro* DNA sequence preferences and confirms that the C-
329 terminus modification does not have significant impact on the specificity of the hairpin
330 polyamides. While based on the pairing rules these results may seem obvious, this experiment
331 was important to confirm that polyamide-conjugates **1** and **2** retain preference for cognate
332 sequences.

333 While PWM-based motifs summarize sequence preferences of DNA-binding molecules, they
334 compress related sequences into a consensus motif, masking the impact of flanking
335 sequences and local microstructure as well as underestimating the affinity spectrum of
336 cognate sites contained within a given DNA-polyamide interactome (DPI). Sequence
337 specificity landscapes (SSLs) can optimize the cognate site motif(s) and thereby uncover
338 major binding motifs to visualize the effects of flanking sequences [40]. To better visualize the
339 full spectrum of DNA binding and compare the individual interactomes of each polyamide, we
340 developed specificity and energy landscapes (SELs) for **1** and **2**, **Fig 2D** [40–42]. SELs
341 present the enriched binding sequences as concentric rings, organized by a “seed motif” in
342 the zero-mismatch ring (central ring) having an exact match for the seed motif, **Fig S4**. The
343 PWM-based motif is used as a seed and the entire DPI displayed in concentric rings as they
344 deviate from the seed motif. Each consecutive ring represents 0, 1, 2, n..., mismatches from
345 the seed motif. SELs plotted for the complete set of enrichment data using a 6-mer seed motif,
346 WGWWCW (**1**) and WGGWCW (**2**), show a clear preference of both polyamides for 6-mer
347 seed motifs (central ring). The dramatic drop-off in affinity for sequences that deviate from the
348 preferred 8-mer site (outer rings), underscores the exquisite sequence specificity of hairpin
349 polyamides, **Fig 2D** and **S5 Fig** [40–42].

350 It is important to note that low-affinity sequences, that are sequentially depleted in SELEX-
351 based approaches, are critical to develop accurate binding site models across the genome
352 [23,37,40,42]. Indeed, we observed both a concentration and selection effect, with each
353 sequential round of SELEX steadily enriching high-affinity sites with a concomitant decrease
354 in correlation with genome-wide binding profiles (**S1 File**). For these reasons, DNA-
355 polyamide interactome from round 1 (with no successive rounds of SELEX) with final
356 concentration of polyamides at 50 nM was used for further analysis.

357

358 **COSMIC-seq to map genome-wide binding profiles**

359 We utilized COSMIC-seq to map the genome-wide binding targets of **1** and **2** in LNCaP nuclei
360 to determine if Py-Im polyamides could maintain their preferred differential binding specificity
361 in a biochemically active complex chromatin environment, **Fig 3A** [24]. Isolated nuclei retain
362 native chromatin states and are widely used to examine chromatin structure and accessibility
363 [50–52]. Briefly, isolated nuclei from human LNCaP cells were treated in biological duplicate
364 with **1** or **2** (0.4 μ M and 4 μ M) at 4 °C for 1 h and cross-linked to DNA by 365 nm UV irradiation.
365 DNA was sheared by sonication, polyamide-DNA complexes were captured by streptavidin-
366 coated magnetic beads, cross-links were reversed, and enriched DNA was sequenced.
367 COSMIC-seq reads were mapped using the Bowtie algorithm and bound peaks were identified

368 using the standard peak-calling algorithms endorsed by the ENCODE consortium, see
369 Methods [24,43,53].

370

371 **Fig 3. Genome-wide DNA binding of polyamides 1 and 2 by COSMIC-seq.** (A) Overview
372 of COSMIC-seq in LNCaP cells, nuclei are treated with polyamides 1 and 2 and cross-linked
373 to DNA with UV irradiation (365 nm). Cross-linked genomic DNA is enriched and analyzed by
374 NGS. (B, C) Heat maps reveal selective enrichment of polyamides 1 and 2. Tag density of
375 each polyamide is shown for the top 1,000 loci for 1 (B) and 2 (C). Data is displayed as
376 sequence read tag density heatmaps (*bottom*) and averaged bar plots (*top*) for the top 1000
377 predicted peaks are mapped on a 10 Kbp window.

378

379 Sequence/Read tag density of the top 1000 identified peaks across replicates of polyamides
380 was compared over a 10 Kbp region centered at bound COSMIC-seq loci, and shown as a
381 heat map, **Fig 3B** and **3C**. Both polyamides 1 and 2 show a strong correlation between bound
382 peaks of replicates while a consistent non-correlation is observed when comparing the top
383 1000 identified peaks of polyamide 1 to those of polyamide 2. A high COSMIC enrichment
384 signal for sites bound by polyamide 1 is observed for replicate treatments (**Fig 3B, left**),
385 however, no enrichment is observed when compared to polyamide 2 (**Fig 3B, right**). A similar
386 negligible enrichment overlap is observed for sites bound by polyamide 2 (**Fig 3C, left**), while
387 a clear enrichment between polyamide 2 replicate treatments is identified (**Fig 3C, right**).
388 Consistent with the *in vitro* binding analysis (**Fig 2C**) polyamide 1 and 2 show poor overlap at
389 bound genomic loci in LNCaP nuclei. These results indicate that differential polyamide
390 sequence specificity is maintained in the context of the native chromatinized genome.

391

392 **Genomic occupancy correlates with binding** 393 **predictions based on summation of sites (SOS)** 394 **model**

395 An overview of the COSMIC (genomic occupancy) and CSI (*in vitro* specificity) pipelines and
396 the comparative analysis of the two datasets is shown in **Fig 4A** (see also Methods). CSI
397 genomescales are generated by assigning binding probability scores across the entire human
398 genome. Unlike standard PWM-based genome annotation approaches, CSI-Genomescales
399 take into account potential moderate-to-low affinity cognate sites [24,40]. Binding intensity is

400 assigned to every 8-bp sequence in the genome based on the CSI data (*in vitro*), and
401 compared to the top 1,000 COSMIC peaks (*in nuclei*) over a 1 Kbp region (**Fig 4A, right**).
402 While the average predicted binding is higher at identified COSMIC loci for polyamides, signal
403 resolution is low when attempting to compare *in vitro* binding at individual COSMIC peaks, **S6**
404 **Fig**. A summation of sites (SOS) scoring model is a more robust method for predicting *in vivo*
405 binding by considering clusters of potential binding sites [24]. Recent studies suggest that
406 genome-wide binding events for natural DNA-binding transcription factors as well as
407 engineered small molecules occur at genomic loci bearing clusters of high-affinity and multiple
408 moderate and weak-affinity binding sites [23]. The CSI *in vitro* binding data of **1** vs **2** clearly
409 illustrate a differential preference for flanking sites and tolerated deviations from seed motifs.
410 These differences could potentially affect the genomic-binding and are known to influence
411 genome-wide distributions of architecturally different polyamides (hairpin versus linear
412 structures). SOS scoring was used to predict binding potential for **1** and **2** across the human
413 genome [23]. By comparing COSMIC signals with predicted binding based on CSI-derived
414 genomescares, we observed that the sum of all *in vitro* determined binding intensities (Z-
415 scores), tiled across an ~420-bp window, most reliably predicted polyamide occupancy at
416 distinct genomic loci (**Fig 4B, S7 Fig, and S2 File**). There is a strong correlation between the
417 top 1000 COSMIC peaks (*in nuclei*) and SOS signals (*in vitro*) for both polyamide **1** and **2**.
418 Notably, there is no correlation observed when comparing *in vitro* and COSMIC data between
419 polyamide **1** and **2**.

420

421 **Fig 4. Polyamide binding across the genome correlates to *in vitro* binding predictions.**
422 (A) COSMIC-scape analysis generates CSI genomescares and SOS heatmaps using the
423 CSI-SELEX and COSMIC-seq data. (B) Data is displayed as averaged bar plots (*top*) and
424 SOS heatmaps (*bottom*) for top 1000 COSMIC peaks mapped on a 10 Kbp region. SOS
425 heatmaps demonstrate that *in vitro* binding of polyamides **1** and **2** is predicted at COSMIC-
426 seq loci while there is no correlation observed between polyamides.

427

428 To access binding differences at individual loci, we selected loci identified by COSMIC on
429 chromosomes 2 and 19 (chr2 and chr19), which were predicted to bind **1** and **2**, respectively,
430 **Fig 5**. By comparing CSI 8-mer enrichment profiles and SOS profiles (**Fig 5A** and **5B**) to the
431 COSMIC tag density (**Fig 5C** and **5D**) for both polyamides, it is evident that polyamides **1** and
432 **2** have distinct, non-overlapping genomic binding preferences. Consistent with sequence-
433 specific binding, polyamide **1** is not found at the CSI-predicted loci on chr19 for polyamide **2**

434 (Fig 5C), and vice versa, **2** is not found at the chr2 predicted loci for **1** (Fig 5D). A similar
435 example of a relationship between loci on chromosome 10 is displayed in S8 Fig. These
436 studies clearly demonstrate a non-correlation between the COSMIC and CSI of **1** vs **2**, and a
437 distinct correlation between *in nuclei* (COSMIC) and *in vitro* (CSI) DNA binding properties.

438

439 **Fig 5. Polyamides 1 and 2 have distinct, non-overlapping genomic binding preferences.**
440 (A, B) Genomescares (*top*) displaying a 1 Kbp region and SOS enrichment plots (*bottom*)
441 displaying a 10 Kbp region for polyamides **1** and **2** at genomic loci of chr2 and chr19,
442 respectively. (C, D) COSMIC tag density data of replicates of **1** and **2** for a 10 Kbp region at
443 same loci.

444

445 DISCUSSION

446 We determined the genome-wide binding events of Py-Im polyamides of similar hairpin
447 architecture that differ at a single position (CH vs N:) on the second aromatic amino acid ring-
448 pair. This single position substitution changes the ring pair from a Py/Py to an Im/Py which
449 alters the binding preference from an A•T or T•A to a G•C base pair, respectively.
450 Comprehensive *in vitro* binding analyses confirm the preferential binding of polyamide **1** to
451 WGWWCW and polyamide **2** to WGGWCW [7]. Both polyamides exhibited high selectivity for
452 their cognate motifs, while exhibiting considerably lower binding affinity for the motif of the
453 other hairpin polyamide. These observations are indicative of polyamide-DNA binding
454 consistent with the established Py-Im pairing rules for both **1** and **2** [54,55]. When comparing
455 binding events within LNCaP nuclei, we see correlation among replicates using the SOS
456 model. Additionally, consistent with the *in vitro* binding analysis, hairpin **1** and **2** show poor
457 overlap in bound genomic loci indicating that innate sequence preferences of these structurally
458 similar hairpin polyamides are maintained in the context of the chromatinized nuclear genome.
459 These results demonstrate that a single position (CH to N:) modification of the aromatic amino
460 acid ring of the polyamide 8-ring structure imparts a significant change in binding preference
461 that is maintained within cellular nuclei.

462 We chose to examine genome-wide binding profiles in intact cell nuclei not only because they
463 present compacted genomic DNA in a chromatinized context but they also circumvent the
464 complexity of cellular uptake of high molecular weight polyamide conjugates. Similarly, state-
465 of-the-art genomic chromatin structure and accessibility studies in a wide range of cell- and

466 tissue-types primarily rely on isolated nuclei because they have been demonstrated to
467 accurately capture native chromatin states in living cells [50–52].

468 Py-Im polyamides have been shown to modulate oncogenic transcription factor signalling and
469 reduce their binding occupancy at select loci in ChIP experiments [14]. A recent study
470 demonstrated that a polyamide targeting dihydrotestosterone (DHT) inducible AR-DNA
471 binding was able to repress 30% of DHT inducible binding events [56]. Importantly, motif
472 analysis of the repressed AR peaks demonstrated that the differential effects on AR-DNA
473 binding events *in vivo* reflects the DNA target sequence binding preference of the hairpin
474 polyamide *in vitro*. Consistent with previous work, we observe a strong correlation between
475 genome-wide SOS scoring using polyamide *in vitro* binding data with the corresponding
476 COSMIC genomic binding data [24].

477 The ability to modify a hairpin Py-Im polyamide while maintaining specificity in a chromatinized
478 environment is an encouraging finding for application as synthetic transcription factors (Syn-
479 TFs). Syn-TFs comprise a modular DNA-binding domain directly fused to or designed to recruit
480 a regulatory domain capable of modulating gene expression when localized to genomic
481 regulatory elements [57,58]. Protein-based artificial TFs have been developed based on zinc
482 finger proteins (ZFPs), transcription activator-like effectors (TALEs), and the clustered
483 regularly interspaced short palindromic repeat (CRISPR)-associated (Cas) system [57,58].
484 However, these methods may be limited by the lack of an efficient delivery mechanism, *in vivo*
485 bioavailability and unknown immunogenic factors [59–61]. The use of small molecule syn-
486 TFs is an attractive non-protein alternative to regulating transcription allowing for more finely
487 tuned control of dosage and timing without the need for complex genomic integration. [43,62–
488 66]. Small molecule solutions have the potential to be used as molecular tools to dissect
489 endogenous gene networks, epigenetic landscapes and as therapeutics to modulate adherent
490 gene expression.

491

492 **Acknowledgements**

493 We thank Laura Vanderploeg for help with the artwork and Mackenzie C. Spurgat for
494 preliminary CSI experiments. Sequencing was performed at the Millard and Muriel Jacobs
495 Genetics and Genomics Laboratory at California Institute of Technology.

496

497 **References**

498 1. Wade WS, Mrksich M, Dervan PB. Design of peptides that bind in the minor groove of

- 499 DNA at 5'-(A,T)G(A,T)C(A,T)-3' sequences by a dimeric side-by-side motif. *J Am*
500 *Chem Soc.* 1992;114(23):8783.
- 501 2. Mrksich M, Wade WS, Dwyer TJ, Geierstanger BH, Wemmer DE, Dervan PB.
502 Antiparallel side-by-side dimeric motif for sequence-specific recognition in the minor
503 groove of DNA by the designed peptide 1-methylimidazole-2-carboxamide netropsin.
504 *Proc Natl Acad Sci U S A.* 1992;89(16):7586.
- 505 3. Trauger JW, Baird EE, Dervan PB. Recognition of DNA by designed ligands at
506 subnanomolar concentrations. *Nature.* 1996;382:559.
- 507 4. White S, Szewczyk JW, Turner JM, Baird EE, Dervan PB. Recognition of the four
508 Watson-Crick base pairs in the DNA minor groove by synthetic ligands. *Nature.*
509 1998;391(6666):468.
- 510 5. Kielkopf CL, White S, Szewczyk JW, Turner JM, Baird EE, Dervan PB, et al. A
511 Structural Basis for Recognition of AT and TA Base Pairs in the Minor Groove of B-
512 DNA. *Science.* 1998;282:111.
- 513 6. Warren CL, Kratochvil NCS, Hauschild KE, Foister S, Brezinski ML, Dervan PB, et al.
514 Defining the sequence-recognition profile of DNA-binding molecules. *Proc Natl Acad*
515 *Sci U S A.* 2006;103(4):867.
- 516 7. Hsu CF, Phillips JW, Trauger JW, Farkas ME, Belitsky JM, Heckel A, et al.
517 Completion of a programmable DNA-binding small molecule library. *Tetrahedron.*
518 2007;63(27):6146.
- 519 8. Minoshima M, Bando T, Sasaki S, Fujimoto J, Sugiyama H. Pyrrole-imidazole hairpin
520 polyamides with high affinity at 5'-CGCG-3' DNA sequence; influence of cytosine
521 methylation on binding. *Nucleic Acids Res.* 2008;36(9):2889.
- 522 9. Puckett JW, Muzikar KA, Tietjen J, Warren CL, Ansari AZ, Dervan PB. Quantitative
523 Microarray Profiling of DNA-Binding Molecules. *J Am Chem Soc.*
524 2007;129(40):12310.
- 525 10. Chenoweth DM, Dervan PB. Allosteric modulation of DNA by small molecules. *Proc*
526 *Natl Acad Sci U S A.* 2009;106(32):13175.
- 527 11. Chenoweth DM, Dervan PB. Structural Basis for Cyclic Py-Im Polyamide Allosteric
528 Inhibition of Nuclear Receptor Binding. *J Am Chem Soc.* 2010;132(41):14521.
- 529 12. Moretti R, Donato LJ, Brezinski ML, Stafford RL, Hoff H, Thorson JS, et al. Targeted
530 Chemical Wedges Reveal the Role of Allosteric DNA Modulation in Protein-DNA
531 Assembly. *ACS Chem Biol.* 2008;3(4):220.
- 532 13. Gottesfeld JM, Neely L, Trauger JW, Baird EE, Dervan PB. Regulation of gene
533 expression by small molecules. *Nature.* 1997;387:202.
- 534 14. Nickols NG, Dervan PB. Suppression of androgen receptor-mediated gene
535 expression by a sequence-specific DNA-binding polyamide. *Proc Natl Acad Sci U S*

- 536 A. 2007;104(25):10418.
- 537 15. Xu L, Wang W, Gotte D, Yang F, Hare AA, Welch TR, et al. RNA polymerase II
538 senses obstruction in the DNA minor groove via a conserved sensor motif. *Proc Natl*
539 *Acad Sci U S A*. 2016;113(44):12426.
- 540 16. Best TP, Edelson BS, Nickols NG, Dervan PB. Nuclear localization of pyrrole–
541 imidazole polyamide–fluorescein conjugates in cell culture. *Proc Natl Acad Sci U S A*.
542 2003;100(21):12063.
- 543 17. Nickols NG, Jacobs CS, Farkas ME, Dervan PB. Improved nuclear localization of
544 DNA-binding polyamides. *Nucleic Acids Res*. 2007;35(2):363.
- 545 18. Yang F, Nickols NG, Li BC, Marinov GK, Said JW, Dervan PB. Antitumor activity of a
546 pyrrole-imidazole polyamide. *Proc Natl Acad Sci U S A*. 2013;110(5):1863.
- 547 19. Martínez TF, Phillips JW, Karanja KK, Polaczek P, Wang C-M, Li BC, et al.
548 Replication stress by Py–Im polyamides induces a non-canonical ATR-dependent
549 checkpoint response. *Nucleic Acids Res*. 2014;42(18):11546–59.
- 550 20. Gottesfeld JM, Melander C, Suto RK, Raviol H, Luger K, Dervan PB. Sequence-
551 specific Recognition of DNA in the Nucleosome by Pyrrole-Imidazole Polyamides. *J*
552 *Mol Biol*. 2001;309(3):615.
- 553 21. Suto RK, Edayathumangalam RS, White CL, Melander C, Gottesfeld JM, Dervan PB,
554 et al. Crystal Structures of Nucleosome Core Particles in Complex with Minor Groove
555 DNA-binding Ligands. *J Mol Biol*. 2003;326(2):371.
- 556 22. Edayathumangalam RS, Weyermann P, Gottesfeld JM, Dervan PB, Luger K.
557 Molecular recognition of the nucleosomal supergroove. *Proc Natl Acad Sci U S A*.
558 2004;101(18):6864.
- 559 23. Erwin GS, Bhimsaria D, Eguchi A, Ansari AZ. Mapping Polyamide–DNA Interactions
560 in Human Cells Reveals a New Design Strategy for Effective Targeting of Genomic
561 Sites. *Angew Chemie Int Ed*. 2014;53(38):10124.
- 562 24. Erwin GS, Grieshop MP, Bhimsaria D, Do TJ, Rodríguez-Martínez JA, Mehta C, et al.
563 Synthetic genome readers target clustered binding sites across diverse chromatin
564 states. *Proc Natl Acad Sci U S A*. 2016;113(47):E7418.
- 565 25. Chandran A, Syed J, Taylor RD, Kashiwazaki G, Sato S, Hashiya K, et al.
566 Deciphering the genomic targets of alkylating polyamide conjugates using high-
567 throughput sequencing. *Nucleic Acids Res*. 2016;44(9):4014.
- 568 26. Anders L, Guenther MG, Qi J, Fan ZP, Marineau JJ, Rahl PB, et al. Genome-wide
569 localization of small molecules. *Nat Biotechnol*. 2014;32(1):92.
- 570 27. Olenyuk BZ, Zhang G-J, Klco JM, Nickols NG, Kaelin JWG, Dervan PB. Inhibition of
571 vascular endothelial growth factor with a sequence-specific hypoxia response element
572 antagonist. *Proc Natl Acad Sci U S A*. 2004;101(48):16768.

- 573 28. Urbach AR, Dervan PB. Toward rules for 1:1 polyamide:DNA recognition. *Proc Natl*
574 *Acad Sci U S A*. 2001;98(8):4343.
- 575 29. Urbach AR, Love JJ, Ross SA, Dervan PB. Structure of a β -Alanine-linked Polyamide
576 Bound to a Full Helical Turn of Purine Tract DNA in the 1:1 Motif. *J Mol Biol*.
577 2002;320(1):55.
- 578 30. Burnett R, Melander C, Puckett JW, Son LS, Wells RD, Dervan PB, et al. DNA
579 sequence-specific polyamides alleviate transcription inhibition associated with long
580 GAA-TTC repeats in Friedreich's ataxia. *Proc Natl Acad Sci U S A*.
581 2006;103(31):11497.
- 582 31. Meier JL, Yu AS, Korf I, Segal DJ, Dervan PB. Guiding the Design of Synthetic DNA-
583 Binding Molecules with Massively Parallel Sequencing. *J Am Chem Soc*.
584 2012;134(42):17814.
- 585 32. Nicholas R. Wurtz, James M. Turner, Eldon E. Baird A, Dervan PB. Fmoc Solid Phase
586 Synthesis of Polyamides Containing Pyrrole and Imidazole Amino Acids. *Org Lett*.
587 2001;3(8):1201.
- 588 33. Sastry SS, Ross BM, P'arraga A. Crosslinking of DNA-binding Proteins to DNA with
589 Psoralen and Psoralen Furan-side Monoadducts. *J Biol Chem*. 1997;272(6):3715.
- 590 34. Puckett JW, Green JT, Dervan PB. Microwave Assisted Synthesis of Py-Im
591 Polyamides. *Org Lett*. 2012;14(11):2774.
- 592 35. Jolma A, Kivioja T, Toivonen J, Cheng L, Wei G, Enge M, et al. Multiplexed massively
593 parallel SELEX for characterization of human transcription factor binding specificities.
594 *Genome Res*. 2010;20(6):861.
- 595 36. Slattery M, Riley T, Liu P, Abe N, Gomez-Alcala P, Dror I, et al. Cofactor binding
596 evokes latent differences in DNA binding specificity between hox proteins. *Cell*. 2011
597 Dec;147(6):1270.
- 598 37. Rodriguez-Martinez JA, Reinke AW, Bhimsaria D, Keating AE, Ansari AZ.
599 Combinatorial bZIP dimers display complex DNA-binding specificity landscapes. *Elife*.
600 2017;6:e19272.
- 601 38. Bailey TL, Williams N, Misleh C, Li WW. MEME: discovering and analyzing DNA and
602 protein sequence motifs. *Nucleic Acids Res*. 2006;34:W369.
- 603 39. Bailey TL, Johnson J, Grant CE, Noble WS. The MEME Suite. *Nucleic Acids Res*.
604 2015;43(W1):W39.
- 605 40. Carlson CD, Warren CL, Hauschild KE, Ozers MS, Qadir N, Bhimsaria D, et al.
606 Specificity landscapes of DNA binding molecules elucidate biological function. *Proc*
607 *Natl Acad Sci U S A*. 2010;107(10):4544.
- 608 41. Tietjen JR, Donato LJ, Bhimisaria D, Ansari AZ. Sequence-Specificity and Energy
609 Landscapes of DNA-Binding Molecules. *Methods Enzymol*. 2011;497:3.

- 610 42. Bhimsaria D, Rodríguez-Martínez JA, Pan J, Roston D, Korkmaz EN, Cui Q, et al.
611 Specificity landscapes unmask submaximal binding site preferences of transcription
612 factors. *Proc Natl Acad Sci U S A*. 2018;115(45):E10586.
- 613 43. Erwin GS, Grieshop MP, Ali A, Qi J, Lawlor M, Kumar D, et al. Synthetic transcription
614 elongation factors license transcription across repressive chromatin. *Science*.
615 2017;358(6370):1617.
- 616 44. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, et al. Simple Combinations
617 of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements
618 Required for Macrophage and B Cell Identities. *Mol Cell*. 2010;38(4):576.
- 619 45. Wade WS, Mrksich M, Dervan PB. Binding Affinities of Synthetic Peptides, Pyridine-2-
620 carboxamidonetropsin and 1-Methylimidazole-2-carboxamidonetropsin, That Form 2:
621 1 Complexes in the Minor Groove of Double-Helical DNA. *Biochemistry*.
622 1993;32:11385.
- 623 46. Kurmis AA, Yang F, Welch TR, Nickols NG, Dervan PB. A Pyrrole-Imidazole
624 Polyamide Is Active against Enzalutamide-Resistant Prostate Cancer. *Cancer Res*.
625 2017;77(9):2207.
- 626 47. Nickols NG, Szablowski JO, Hargrove AE, Li BC, Raskatov JA, Dervan PB. Activity of
627 a Py-Im Polyamide Targeted to the Estrogen Response Element. *Mol Cancer Ther*.
628 2013;12(5):675.
- 629 48. Kashiwazaki G, Chandran A, Asamitsu S, Kawase T, Kawamoto Y, Sawatani Y, et al.
630 Comparative Analysis of DNA-Binding Selectivity of Hairpin and Cyclic Pyrrole-
631 Imidazole Polyamides Based on Next-Generation Sequencing. *ChemBioChem*.
632 2016;17(18):1752.
- 633 49. Tuerk C, Gold L. Systematic Evolution of Ligands by Exponential Enrichment: RNA
634 Ligands to Bacteriophage T4 DNA Polymerase. *Science*. 1990;246(4968):505.
- 635 50. Boyle AP, Davis S, Shulha HP, Meltzer P, Margulies EH, Weng Z, et al. High-
636 Resolution Mapping and Characterization of Open Chromatin across the Genome.
637 *Cell*. 2008;132(2):311.
- 638 51. Dunham I, Kundaje A, Aldred SF, Collins PJ, Davis CA, Doyle F, et al. An integrated
639 encyclopedia of DNA elements in the human genome. *Nature*. 2012;489(7414):57.
- 640 52. Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of
641 native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-
642 binding proteins and nucleosome position. *Nat Methods*. 2013;10(12):1213.
- 643 53. Landt SG, Marinov GK, Kundaje A, Kheradpour P, Pauli F, Batzoglou S, et al. ChIP-
644 seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome*
645 *Res*. 2012;22(9):1813.
- 646 54. Dervan PB, Kurmis AA, Finn PB. Chapter 12. Molecular Recognition of DNA by Py-

- 647 Im Polyamides: From Discovery to Oncology. DNA-targeting Mol as Ther Agents.
648 2018;(7):298.
- 649 55. White S, Baird EE, Dervan PB. On the pairing rules for recognition in the minor
650 groove of DNA by pyrrole-imidazole polyamides. Chem Biol. 1997;4(8):569.
- 651 56. Kurmis AA, Dervan PB. Sequence specific suppression of androgen receptor–DNA
652 binding in vivo by a Py-Im polyamide. Nucleic Acids Res. 2019;47(8):3828.
- 653 57. Heiderscheidt EA, Eguchi A, Spurgat MC, Ansari AZ. Reprogramming cell fate with
654 artificial transcription factors. FEBS Lett. 2018;592(6):888.
- 655 58. Black JB, Gersbach CA. Synthetic transcription factors for cell fate reprogramming.
656 Curr Opin Genet Dev. 2018;52:13.
- 657 59. Kadi AM, Gersbach CA. Engineering synthetic TALE and CRISPR/Cas9
658 transcription factors for regulating gene expression. Methods. 2014;69(2):188.
- 659 60. Crudele JM, Chamberlain JS. Cas9 immunity creates challenges for CRISPR gene
660 editing therapies. Nat Commun. 2018;9(1):3497.
- 661 61. Charlesworth CT, Deshpande PS, Dever DP, Camarena J, Lemgart VT, Cromer MK,
662 et al. Identification of preexisting adaptive immunity to Cas9 proteins in humans. Nat
663 Med. 2019;25(2):249.
- 664 62. Mapp AK, Ansari AZ, Ptashne M, Dervan PB. Activation of gene expression by small
665 molecule transcription factors. Proc Natl Acad Sci U S A. 2000;97(8):3930.
- 666 63. Ansari AZ, Mapp AK, Nguyen DH, Dervan PB, Ptashne M. Towards a minimal motif
667 for artificial transcriptional activators. Chem Biol. 2001;8(6):583.
- 668 64. Arora PS, Ansari AZ, Best TP, Ptashne M, Dervan PB. Design of Artificial
669 Transcriptional Activators with Rigid Poly-l-proline Linkers. J Am Chem Soc.
670 2002;124:13067.
- 671 65. Kwon Y, Arndt H-D, Mao Q, Choi Y, Kawazoe Y, Dervan PB, et al. Small Molecule
672 Transcription Factor Mimic. J Am Chem Soc. 2004;126:15940.
- 673 66. Kawamoto Y, Bando T, Sugiyama H. Sequence-specific DNA binding Pyrrole–
674 imidazole polyamides and their applications. Bioorg Med Chem. 2018;26(8):1393–
675 411.

676

677 Supporting information

- 678 **S1 Fig. Chemical Synthesis and characterization of Py-Im polyamide conjugates. (A)**
679 Solid phase synthetic scheme for the synthesis of Py-Im polyamides **1A** and **2A**, **a**) Boc-Py-
680 OBt, *i*-Pr₂NEt, DMF, μ W (80 °C, 3 h); **b**) 80:1:19 TFA:triethylsilane:CH₂Cl₂, 5 min, RT; **c**) Boc-
681 Py-OH, PyBOP *i*-Pr₂NEt, DMF, μ W (60 °C, 5 min); **d**) 9:2:1 DMF:Ac₂O:*i*-Pr₂NEt, 30 min, RT;
682 **e**) repeat (1x) steps b - d; **f**) 80:1:19 TFA:triethylsilane:CH₂Cl₂, 5 min, RT; **g**) Boc-Im-OH,

683 PyBOP, *i*-Pr₂NEt, DMF, μ W (60 °C, 5 min); **h**) 9:2:1 DMF:Ac₂O:*i*-Pr₂NEt, 30 min, RT; **i**) 80:1:19
684 TFA:triethylsilane:CH₂Cl₂, 25 min, RT; **j**) Fmoc-D-Dab(Boc)-OH, PyBOP, *i*-Pr₂NEt, DMF, μ W
685 (60 °C, 25 min); **k**) 9:2:1 DMF:Ac₂O:*i*-Pr₂NEt, 30 min, RT; **l**) repeat (2x) steps b - d; **m**) 80:1:19
686 TFA:triethylsilane:CH₂Cl₂, 5 min, RT; **n**) **5**, PyBOP, *i*-Pr₂NEt, DMF, μ W (60 °C, 5 min); **o**) 20%
687 piperidine, DMF, 30 min, RT; **p**) 9:2:1, DMF:Ac₂O:*i*-Pr₂NEt 30 min, RT; **q**) neat 3,3'-Diamino-
688 N-methyldipropylamine, μ W (60 °C, 10 min); **(B)** Synthesis of the psoralen–biotin-acid moiety
689 **3**, **a**) Fmoc-PEG₂-OH, *i*-Pr₂NEt, DCM; **b**) 20% piperidine, DMF; **c**) Biotin-Lys(Fmoc)-OH,
690 HATU, HOAt, *i*-Pr₂NEt, 3:1 DMSO:DMF; **d**) 20% piperidine, DMF; **e**) Fmoc-PEG₂-OH, HATU,
691 HOAt, *i*-Pr₂NEt, DMF; **f**) 20% piperidine, DMF; **g**) SPB (NHS-psoralen), *i*-Pr₂NEt, DMF; **h**) 95%
692 TFA, 2.5% H₂O, 2.5% *i*-Pr₃SiH; and **(C)** peptide coupling of Py-Im polyamides **1A** and **2A** with
693 **3**; **(C)** Analytical HPLC traces of **1**, **2**, and **3**; **(D)** Characterization of compounds by MALDI-
694 TOF.

695

696 **S2 Fig. Position weight matrix (PWM) motif logo representation obtained from the top**
697 **50 sequences via CSI by SELEX-seq for polyamides conjugates.** PWMs for three
698 replicates of **1** (*left*) and **2** (*right*) at two concentrations (5 nM and 50 nM) and three enrichment
699 rounds (1, 2 and 3). The PWMs are derived using MEME software with the corresponding e-
700 value indicated.

701

702 **S3 Fig. Specificity and energy landscapes (SELs) display the comprehensive binding**
703 **preferences of polyamides based on a seed motif.** The height of each peak corresponds
704 to CSI enrichment for a given sequence. **(A)** SEL for **1** with seed motif WGWWCW (where W
705 = A or T). **(B)** Top view of the SEL in **A**. **(C)** SEL for **2** with seed motif WGGWCW (where W
706 = A or T). **(D)** Top view of the SEL in **C**. **(E)** SELs consists of concentric rings with sequences
707 in the 0 mismatch ring (central ring) having an exact match to the seed motif. Moving outwards,
708 the 1 mismatch ring contains all sequences that differ from the seed motif at any one position
709 (or a Hamming distance of one). In each ring, sequences are arranged clockwise by position
710 of the mismatch, then alphabetically by the sequence. The 1 mismatch ring begins with
711 mismatches at the first position of the motif and ends with mismatches at the last position of
712 the motif.

713

714 **S4 Fig. Specificity and energy landscapes (SEL) representation for all k-mer binding**
715 **enrichment obtained via CSI by SELEX-seq for polyamide conjugates.** SELs for three
716 replicates of **1** (*left*) and **2** (*right*) at two concentrations (5 nM and 50 nM) and three enrichment
717 rounds (1, 2 and 3).

718

719 **S5 Fig. Scatter plots and correlation coefficients denoting replicability of CSI replicates**
720 **for polyamides conjugates.** Scatter plots for CSI enrichment of **1 (A)** and **2 (B)** at two
721 concentrations (5 nM and 50 nM) and three enrichment rounds (1, 2 and 3).

722

723 **S6 Fig. COSMIC loci compared to CSI genomescales.** Data is displayed as averaged bar
724 plots (*top*) and heatmap of genomescales (*bottom*) for top 1000 COSMIC peaks mapped on
725 a 1 Kbp region. CSI data from enrichment round 1 at 50 nM for **1 (left)** and **2 (right)** was used
726 for genomescale generation.

727

728 **S7 Fig. COSMIC-seq tag density data plotted for the top 1000 SOS predicted sites shows**
729 **COSMIC binding is found at the predicted genomic sites.** Heatmaps with tag density for
730 COSMIC replicates of **1 (A)** and **2 (B)** are mapped for the top 1000 SOS predicted genomic
731 peaks using a 10 Kbp window. CSI data from enrichment round 1 at 50 nM for **1 (left)** and **2**
732 (*right*) was used for SOS prediction.

733

734 **S8 Fig. Polyamides 1 and 2 have distinct, non-overlapping genomic binding**
735 **preferences.** (A) Genomescales (*top*) displaying a 1 Kbp region and SOS enrichment plots
736 (*bottom*) displaying a 10 Kbp region for polyamides **1** and **2** at genomic loci of chr10. (B)
737 COSMIC tag density data of replicates of **1** and **2** for a 10 Kbp region at same loci.

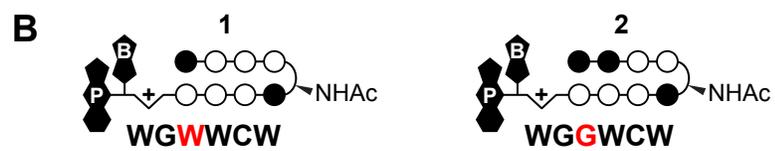
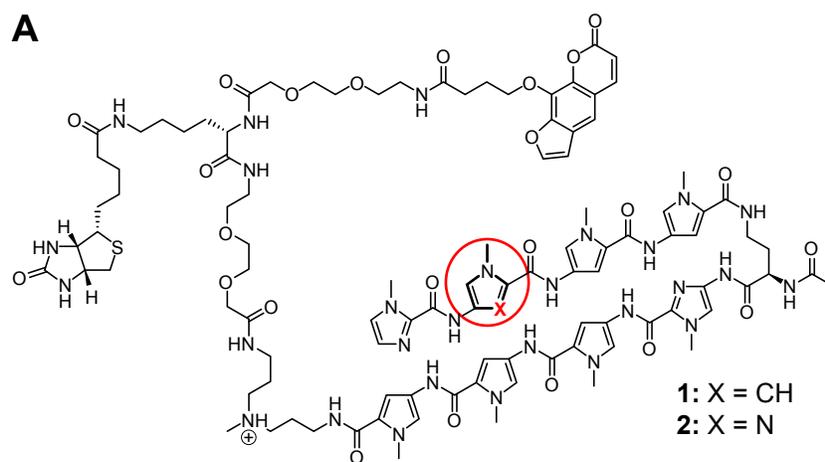
738

739 **S1 File. SOS and genomescales.** Heatmaps for SOS and genomescale data for enrichment
740 round 1 at 50 nM for polyamide **1** and **2**. Heatmaps are plotted for the top 1000 COSMIC
741 peaks of COSMIC replicates of polyamides **1** and **2** on a 10 Kbp window for SOS and 1 Kbp
742 for genomescales.

743

744 **S2 File. COSMIC-seq tag density data.** Tag density heatmaps for polyamide **1** and **2**
745 replicates are mapped for the top 1000 SOS predicted genomic peaks using a 10 Kbp window.
746 CSI data from enrichment round 1 at 50 nM for **1** and **2** was used for SOS prediction.

Figure 1



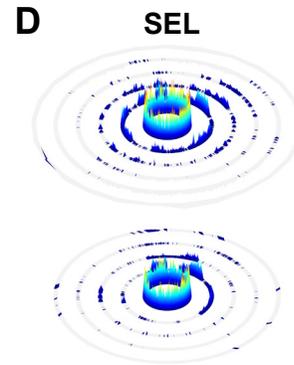
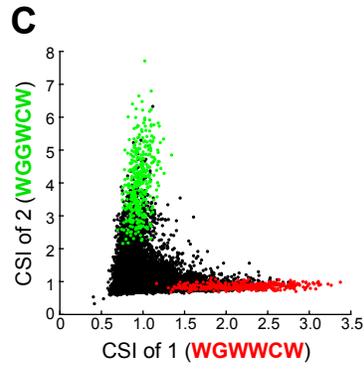
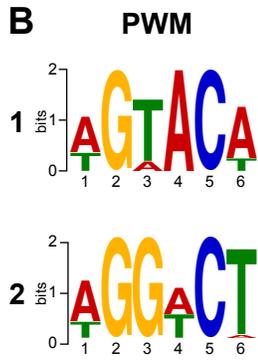
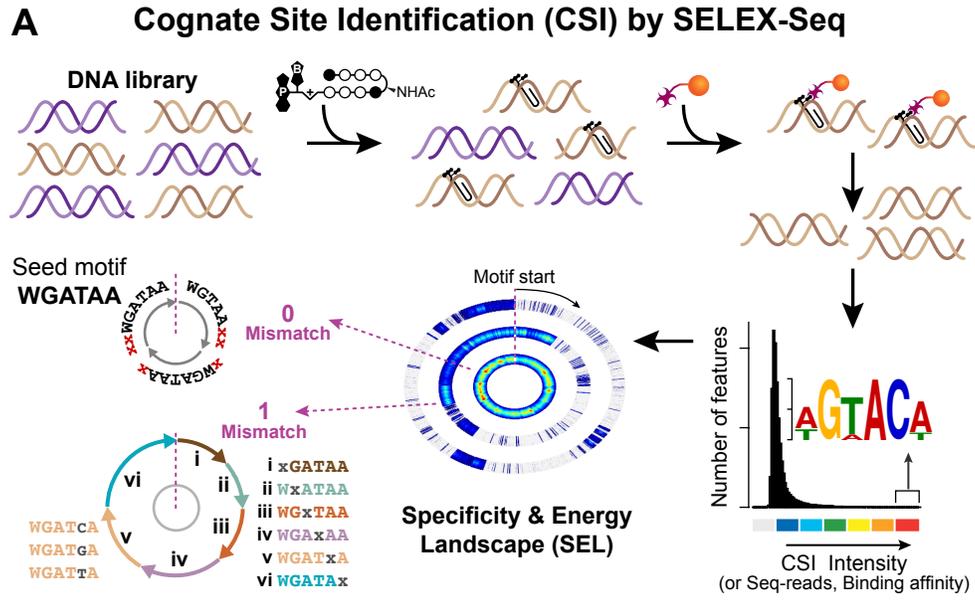


Figure 3

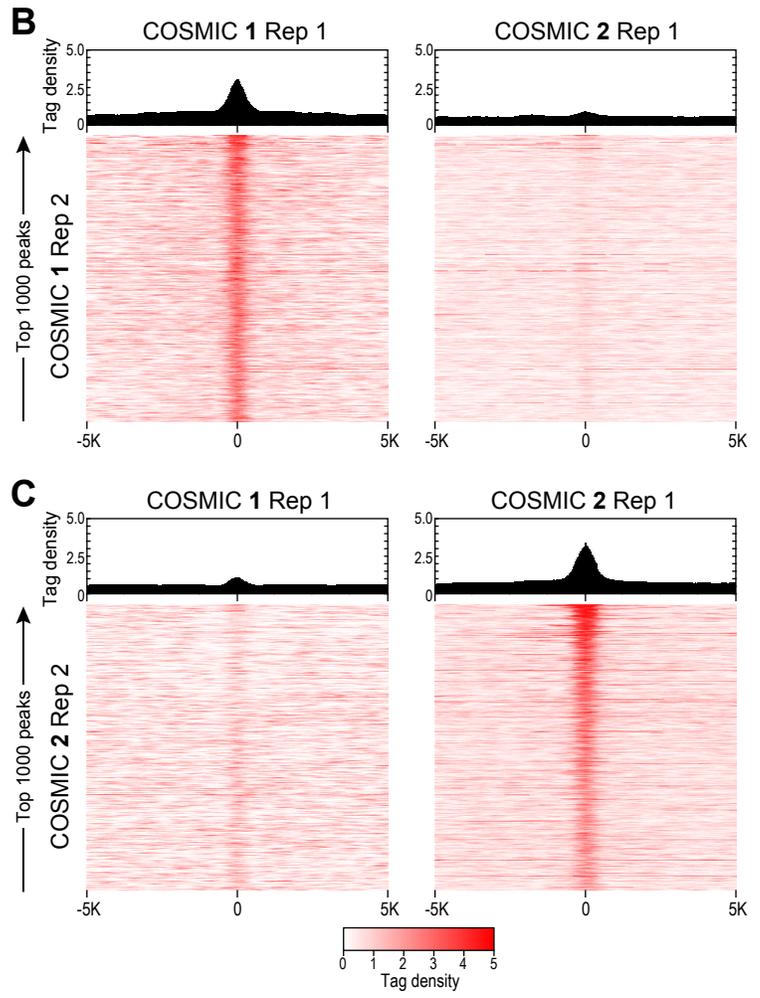
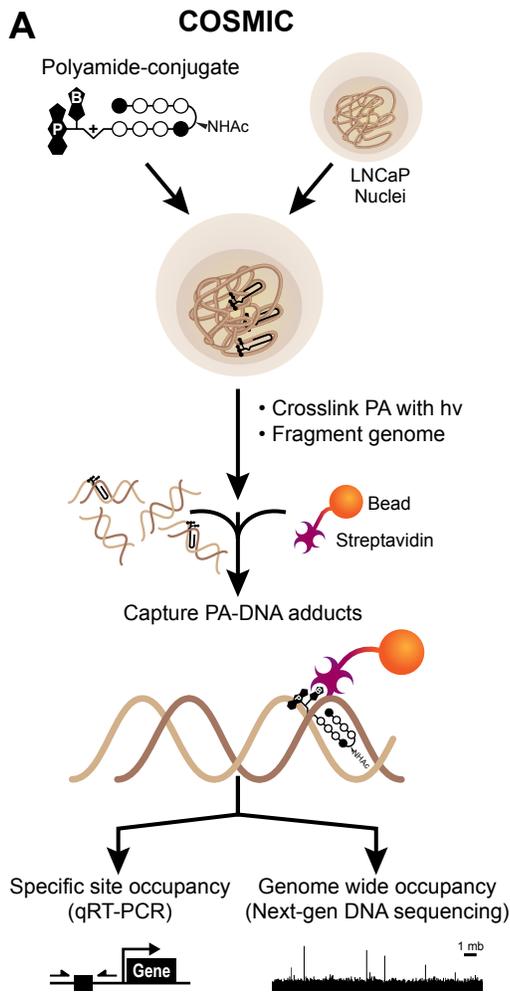


Figure 4

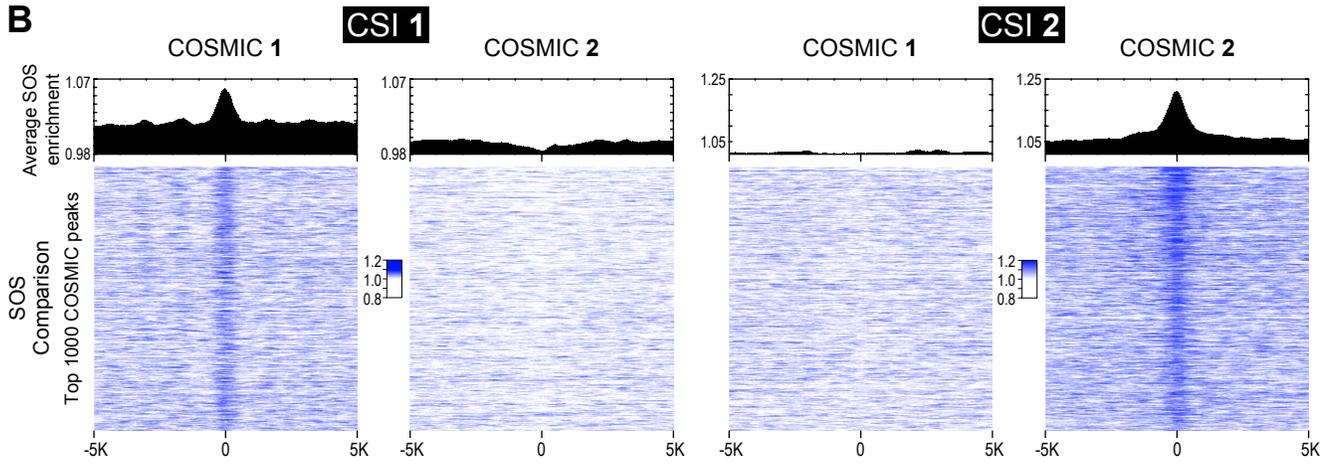
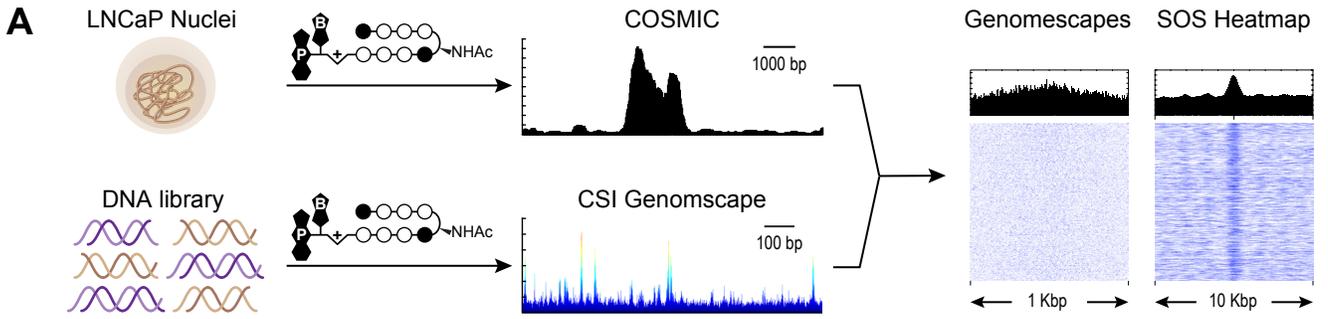


Figure 5

