
COMPOUND-SPECIFIC DNA ADDUCT PROFILING WITH NANOPORE SEQUENCING AND IONSTATS

Yrjö Koski^{1,2}, Divyesh Patel^{2,3}, Natalia Kakko von Koch⁴, Paula Jouhten⁴, Lauri Aaltonen^{2,5,6},
Kimmo Palin^{2,5,6,*}, Biswajyoti Sahu^{2,3,7,*}, and Esa Pitkänen^{1,2,6,*}

¹Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland

²Applied Tumor Genomics Research Program, Faculty of Medicine, University of Helsinki, Helsinki, Finland

³The Norwegian Centre for Molecular Biosciences and Medicine (NCMBM), University of Oslo, Oslo, Norway

⁴Department of Bioproducts and Biosystems, School of Chemical Engineering, Aalto University, Espoo, Finland

⁵Department of Medical and Clinical Genetics, University of Helsinki, Helsinki, Finland

⁶iCAN Digital Precision Cancer Medicine Flagship, Helsinki, Finland

⁷Department of Cancer Genetics, Institute for Cancer Research, Oslo, Norway

*Corresponding authors. Email: kimmo.palin@helsinki.fi, biswajyoti.sahu@helsinki.fi,
esa.pitkanen@helsinki.fi

ABSTRACT

1 Covalently bound DNA adducts are mutation precursors that contribute to aging and diseases such as
2 cancer. Accurate detection of adducts in the genome will shed light on tumorigenesis. Commonly
3 used detection methods are unable to pinpoint the exact genomic locations of adducts. Long-
4 read nanopore sequencing has the potential to accurately detect multiple types of DNA adducts
5 at single-nucleotide precision. In this study, we developed a novel statistical toolkit, IonStats, to
6 profile DNA adducts in nanopore sequencing data. With IonStats, we investigated the effects of
7 four adduct-inducing genotoxic compounds on nanopore sequencing, and found both shared and
8 compound-specific perturbations in base quality scores, ionic current profiles, and translocation
9 dynamics. Notably, aristolochic acid II and melphalan treatments profoundly altered nanopore
10 readouts and led to substantial sequence-specific read interruptions. Our study shows that nanopore
11 sequencing can be effectively employed to detect and characterize DNA adducts, paving the way for
12 high-resolution, high-throughput profiling of DNA damage and the exposome.

13 **Keywords** Nanopore sequencing · DNA adducts · DNA adduct detection · genotoxicity · IonStats

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

14 **1 Introduction**

15 DNA adducts are covalently bound lesions caused by exogenous and endogenous genotoxic compounds that react
16 with DNA directly or after metabolic activation [1]. Enzymatic DNA repair mechanisms remove most adducts, but
17 a small fraction may escape the repair [2]. Adducts can interfere with DNA transcription and replication, and give
18 rise to mutations that may lead to cancer [3]. Several compounds have been implicated in cancer initiation, including
19 aristolochic acids (AAs), polycyclic aromatic hydrocarbons (PAHs), and aflatoxin B₁ (AFB₁), linked to urological, lung,
20 and liver cancers, respectively [4, 5, 6, 7]. Beyond cancer, DNA adducts serve as molecular imprints of environmental
21 exposures, contributing to accelerated biological aging and increased mortality [8]. Detecting and characterizing DNA
22 adducts is thus crucial for understanding how exposure to DNA damage can lead to mutations and cellular dysfunction
23 underlying both carcinogenesis and the progressive decline associated with aging.

24 There are two distinct classes of DNA adducts: monoadducts and crosslinks [9]. Monoadducts are formed through
25 chemical modifications by small molecules on individual nucleotides. The primary targets of monoadduct binding
26 depend on the reactive species, nucleophilicity of DNA sites, steric factors, and DNA sequence context. For small
27 alkylating agents, the primary binding sites are the ring nitrogens N3 and N7, which are most nucleophilic in purine
28 bases. For bulky aromatic compounds, such as AAs, the primary binding sites are exocyclic amino groups [10].
29 Molecules with multiple reacting groups can form DNA-DNA crosslinks, where two nucleotides are bound by the
30 adduct [11]. Crosslinks can be either inter- or intrastrand, where the bound nucleotides are on different strands or the
31 same strand, respectively [12, 13].

32 The most prominent methods for DNA adduct detection include mass spectrometry- and amplicon-based meth-
33 ods [14]. High-resolution mass spectrometry enables comprehensive profiling of adduct exposure, detecting most DNA
34 monoadducts that are present in the sample [15, 16]. However, these methods require the DNA to be broken down into
35 individual nucleotides, losing the positional information of the adducts in the genome [14]. Amplicon-based methods
36 solve this issue by recognizing the adducts chemically or enzymatically by using DNA repair enzymes to mark and
37 cleave DNA adducts [17, 18]. However, amplicon-based methods are unable to distinguish between different types of
38 adducts [14].

39 Nanopore sequencing is an emerging technology with the potential to differentiate between adduct types at single-
40 nucleotide resolution while preserving site-specificity [19, 20]. Nanopore sequencing measures the ionic current as
41 a single-stranded DNA fragment translocates through a nanopore, which is embedded in an electrically insulating
42 membrane. The nucleobases that reside inside the nanopore at any given time alter the ionic current based on their
43 physical structure. Hence, analysis of the ionic current can potentially reveal both the DNA sequence and modifications
44 such as adducts [21, 22].

45 Previous studies have focused on adducts only in specific known sequence contexts [19, 20]. In this study, we show that
46 nanopore sequencing can be used to characterize adducts in arbitrary sequence contexts for multiple adduct-inducing
47 molecules. We implemented a novel software package called IonStats for distilling nanopore sequencing data into

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

48 multiple statistics revealing deviations from the unadducted state. We show that these statistics are informative for adduct
49 formation by creating and analyzing a novel nanopore sequencing dataset of whole genome amplified *Saccharomyces*
50 *cerevisiae* DNA with four adduct-inducing genotoxic compounds: aristolochic acid II (AAlI), cisplatin, melphalan,
51 and mitomycin C (MMC). These compounds vary in size and chemistry, and bind to various nucleotides in DNA.
52 We demonstrate distinct treatment effects by comparing the sequencing statistic distributions between treated and
53 untreated samples. We also observed that some treatments increase read interruptions, where DNA fragments are not
54 fully sequenced, likely due to bulky adducts or crosslinks bound at specific sequence contexts. IonStats is available
55 open source at GitHub (<https://github.com/ykoski/ionstats>).

56 **2 Results**

57 **2.1 IonStats: a comprehensive toolkit for capturing adduct effects in nanopore sequencing data**

58 Nanopore sequencing generates large volumes of ionic current signal data, posing substantial challenges for downstream
59 analyses (Fig. 1A,B). To distill informative features for adduct detection, we developed a computational workflow
60 called IonStats (Methods). IonStats collects and summarizes diverse statistics from nanopore signal and basecalled
61 sequencing data by k -mers or reference genome positions (Methods; Fig. 1C). Comparison of the statistics of treated
62 samples with those of untreated samples allows the characterization of treatment effects. In addition, IonStats collects
63 read-level statistics including read counts and lengths, and base quality scores. In our study, we used $k = 9$, or the
64 number of nucleotides that reside within the nanopore (ONT pore version R10.4.1) at any given time [23]. This choice
65 also ensured high k -mer coverage in our experiments.

66 For each k -mer or reference position, IonStats collects nanopore signal means and standard deviations, base quality
67 scores, and dwell times. The latter indicates the number of measurements obtained for each nucleotide, correlating with
68 the time the nucleotide resided in the nanopore sensing region (Fig. 1A,C). In addition, IonStats computes motor-protein
69 adjusted dwell time (ADT). This statistic reflects the duration that each nucleotide resides at the motor protein, rather
70 than within the sensing region. ADT lets us study the effects that adducts might have when interacting with the motor
71 protein, in addition to in-pore effects. IonStats also provides a novel windowed error score called Average Expected
72 Absolute Deviation (AEAD; Methods). AEAD captures how much the signal means or standard deviations differ from
73 expected values within a k -mer window. Finally, adducts may cause failed translocations, which IonStats can detect by
74 identifying interrupted reads through missing rear barcode sequences.

75 IonStats provides distributional tests to detect differences between sample groups, typically between treated and control
76 samples. The available tests include the Kolmogorov-Smirnov (KS) test and the Cramér von Mises (CvM) criterion.
77 Compared to the more common KS test, CvM is more sensitive to differences in tails of distributions. Together, these
78 complementary tests enable robust detection of treatment-induced perturbations in nanopore data.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

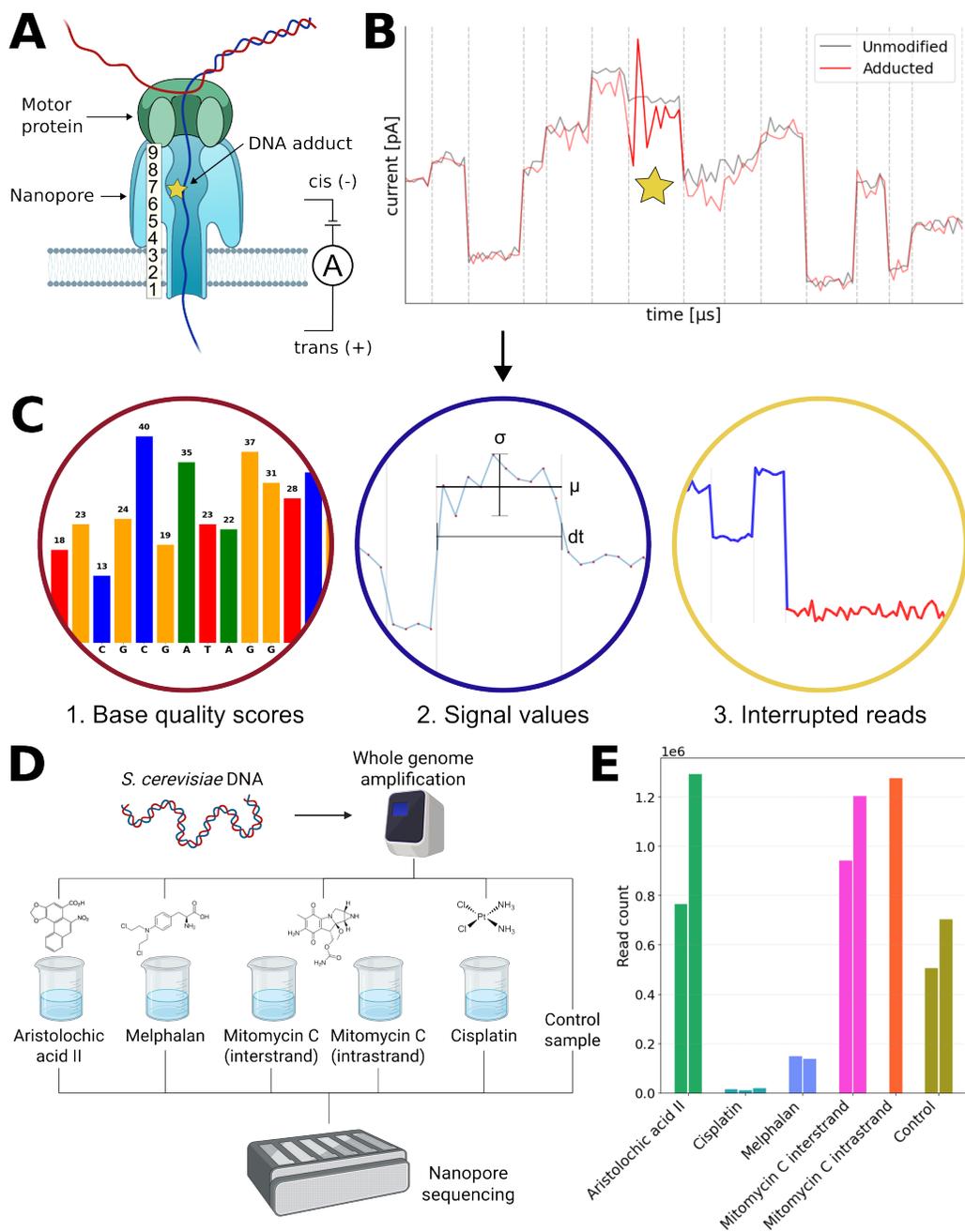


Figure 1: Key concepts and study design. **A**: Visualization of the nanopore sequencing process. Positions of the nanopore are numbered. **B**: An illustration of the effect that a DNA adduct can have on the ionic current. **C**: Overview of the data analysis approaches that IonStats uses to distill nanopore sequencing data into multiple statistics to compare treated and untreated control samples. **D**: Visualization of the experimental setup of this study. **E**: Number of reads per sample replicate obtained in the experiment.

79 2.2 Adduct-induced DNA damage reduces nanopore read length and quality

80 We first examined whether the effects of adduct treatments can be observed in read-level nanopore sequencing statistics
81 after treating *S. cerevisiae* DNA with four genotoxic, adduct-inducing compounds (Fig. 1D; Methods). These statistics

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

82 included the number of reads sequenced (Fig. 1E), read lengths and mean base quality scores (Fig. 2), and rear barcode
83 quality scores (Suppl. Fig. 1, Suppl. Table 1). AAI, melphalan, and cisplatin samples showed lower mean quality
84 scores (16.6, 17.8, and 5.23, respectively; Suppl. Table 2), whereas in mitomycin samples, scores were higher (MMC
85 interstrand 19.9, MMC intrastrand 19.9) than in untreated controls (19.6, $p < 10^{-16}$ for all compounds) (Fig. 2). Both
86 melphalan and cisplatin samples had shorter reads (median lengths: 1833 and 3658, respectively) compared to other
87 samples (AAI 6744, MMC interstrand 7430, MMC intrastrand 7331, control 7061; Suppl. Table 2). The replicates
88 generally correlated very well, except for AAI, where the two replicates differed by mean quality score (means 15.8
89 and 17.1; $p < 10^{-16}$).

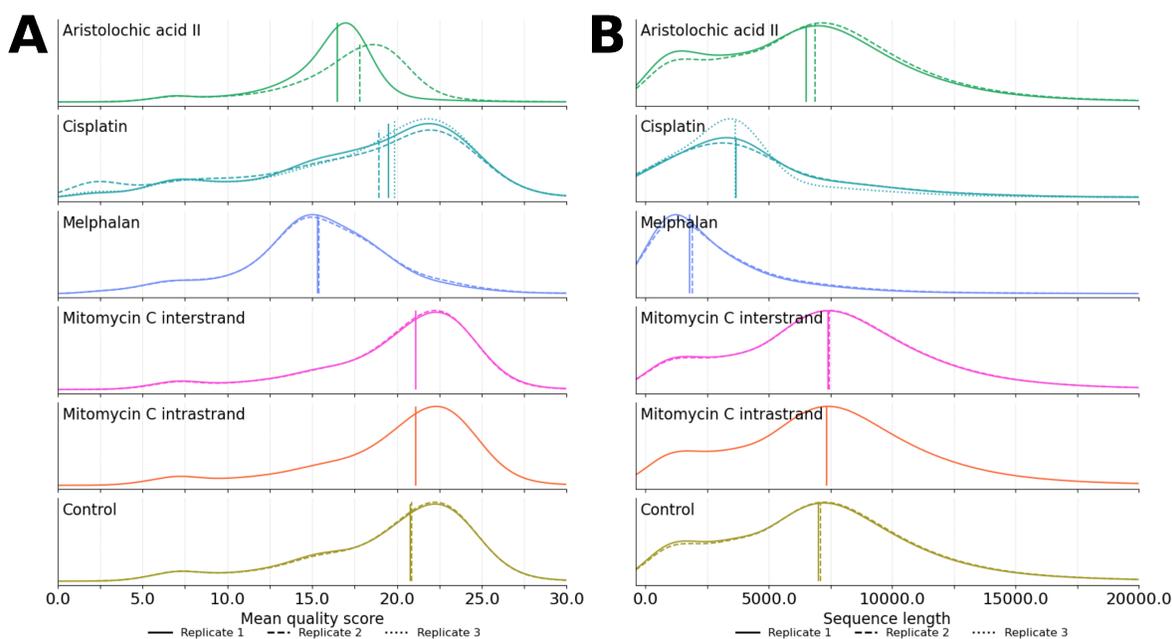


Figure 2: Overview of the sequencing statistics in the nanopore sequencing dataset. **A:** Read mean quality scores by sample. **B:** Read lengths by sample. The vertical line shows the median value of each replicate.

90 2.3 Treatment effects are observed in tails of nanopore signal-level statistics

91 We then focused on the signal-level statistics obtained with IonStats. Overall, for most treatments, the distributions of
92 these k -mer-level statistics were consistent with those observed in controls (Suppl. Table 3 & 4, Suppl. Fig. 2). However,
93 for AAI, melphalan, and cisplatin, the variability of many statistics differed substantially from controls (Suppl. Table
94 3 & 5, Suppl. Fig. 3). This led us to believe that the treatment effects are primarily observed at the tails of k -mer
95 value distributions. Consistent with this hypothesis, the extreme quantiles (1st and 99th percentiles) of signal statistics
96 revealed pronounced shifts (Fig. 3, Suppl. Table 5), while their means showed only minor differences (Suppl. Table 3,
97 4). This supported the notion that adduct formation affects only a subset of DNA nucleotides, producing a mixture of
98 modified and unmodified k -mers (Fig. 3 A, B). Due to insufficient k -mer coverage, cisplatin samples were excluded
99 from these analyses (Suppl. Table 1).

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

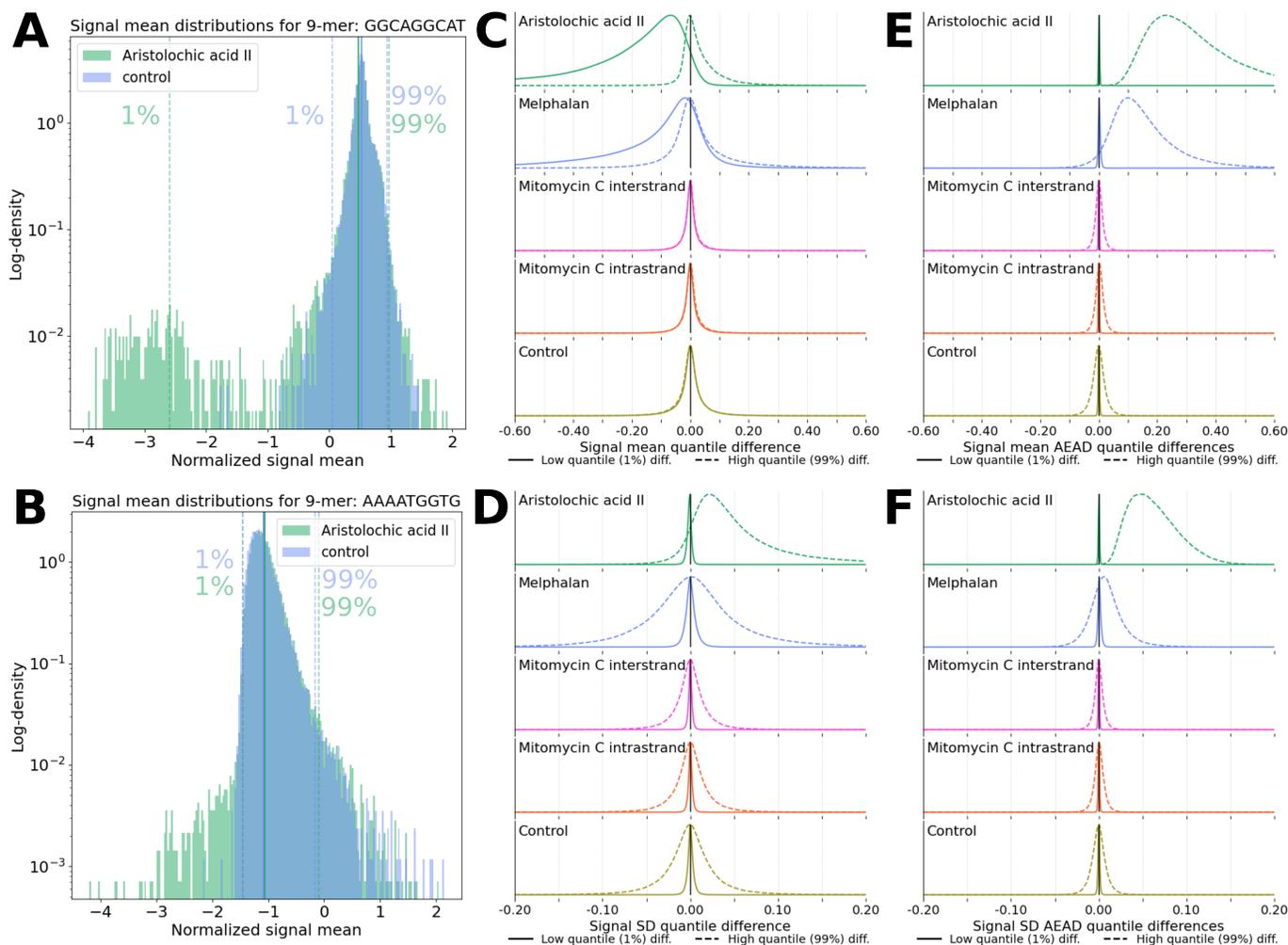


Figure 3: **A:** Nanopore signal mean distribution for the 9-mer GGCAGGCAT, where 1% quantile outliers distinguish the AAI-treated sample from controls. Dashed lines and annotations indicate the 1% and 99% quantiles; solid lines indicate the distribution means. **B:** An example of a 9-mer (AAAATGGTG) where the signal mean distribution is non-normal but highly similar between AAI-treated and control samples, with no substantial differences in the extreme quantiles. **C-F:** Quantile differences for signal mean (**C**), standard deviation (**D**), mean AEAD (**E**), and standard deviation AEAD (**F**). For treatments, differences for each k -mer are shown between treatment and controls; for controls, between-replicate differences are shown.

100 Among all treatments, we observed the strongest effects in AAI-treated samples, with melphalan samples showing
 101 similar but weaker effects (Fig. 3C-F; Suppl. Table 4-7). Specifically, the lowest quantile signal means for most k -mers
 102 in AAI samples were lower than those in controls (Fig. 3C). This is consistent with bulky DNA adducts impeding ion
 103 flow through the nanopore, thus lowering the measured current. We also found that many k -mers displayed markedly
 104 higher variability of signal levels, particularly in AAI samples (Fig. 3D). Furthermore, aggregating signal over multiple
 105 consecutive k -mers, our AEAD score displayed more pronounced alterations in AAI and melphalan samples than the
 106 other statistics (Fig. 3E,F). While the effects of adduct treatments were evident in both AAI and melphalan samples
 107 (Suppl. Table 4), signal-level statistics did not differentiate mitomycin C samples from controls.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

108 2.4 Motif discovery reveals sequences associating with treatment effects

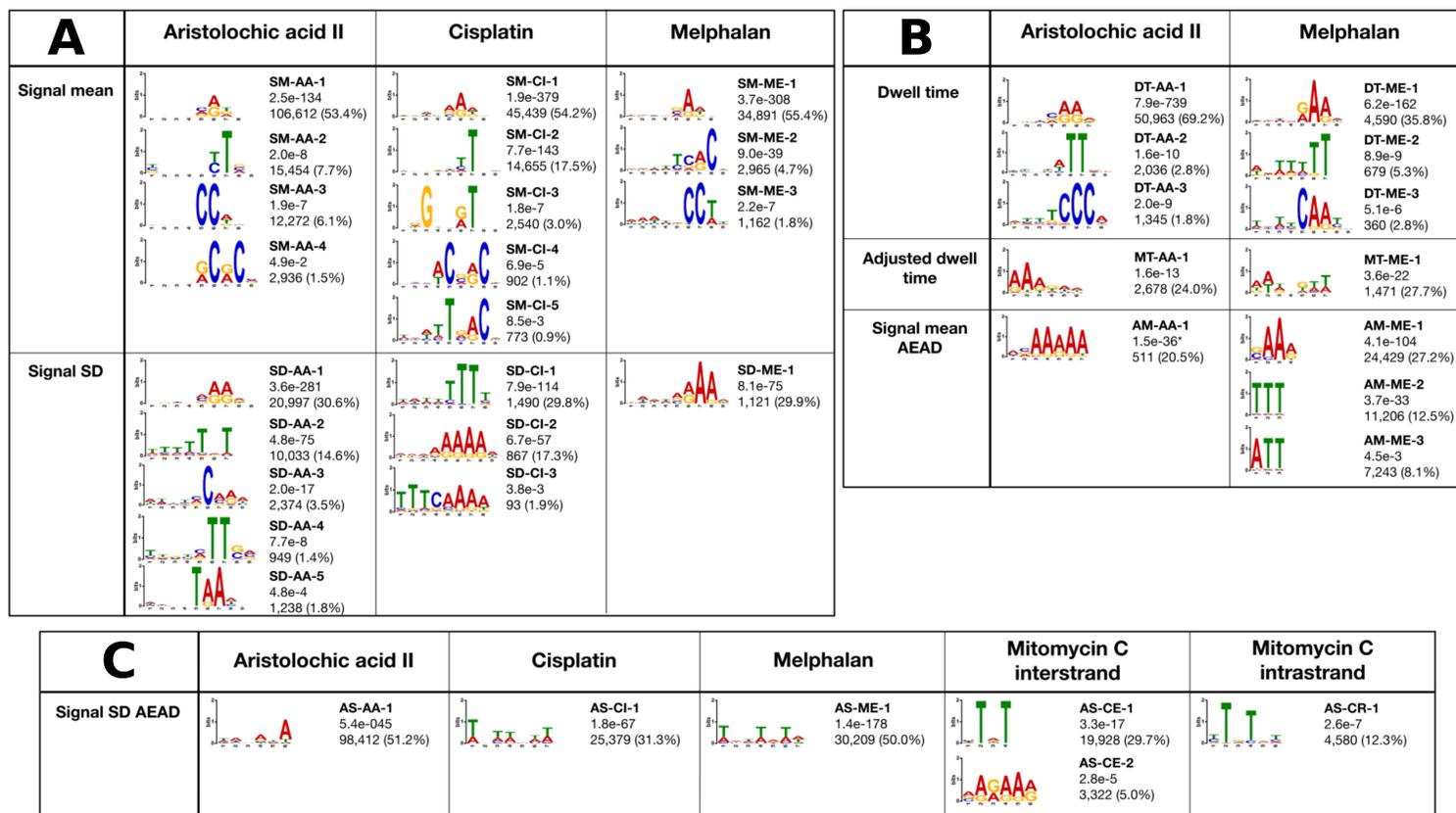


Figure 4: All significant, high-quality motifs found with the STREME analysis. A motif had to occur on both strands and have an E-value smaller than 0.05 to be included. For each motif, its identifier (top), E-value (middle), and the number and fraction of motif matches in the significant 9-mers (bottom) are shown. Significant motifs shown for **A**: signal mean and standard deviation, **B**: dwell times and signal mean AEAD, and **C**: signal standard deviation AEAD. In **B**, * refers to a reduced k -mer set (Methods).

109 Next, we asked whether treatment-driven effects might be linked to recurring DNA sequence contexts. To explore
 110 this, we conducted sequence motif discovery on k -mers associated with differences in sequencing statistics between
 111 the treated and control samples, as determined by the CvM criterion (Methods, Suppl. Fig. 4). This analysis revealed
 112 multiple significant motifs ($E < 0.05$), some specific to individual treatments and others shared across treatments
 113 (Fig. 4). No significant motifs were linked to base quality scores in any of the treatments.

114 We identified motifs in AAIL, cisplatin, and melphalan samples to associate with the signal mean level and variability
 115 (Fig. 4A). The most prominent motifs (SM-AA-1, SM-CI-1, SM-ME-1; Fig. 4) all showed purine preference at upper
 116 pore positions (*i.e.*, closest to the motor protein; Fig. 1A).

117 We observed pyrimidine-rich motifs at upper pore positions in AAIL and cisplatin samples (SM-AA-2, SM-CI-2).
 118 Additionally, weaker motifs containing pairs of cytosines, either directly adjacent or separated by one to two nucleotides,
 119 were also identified (SM-AA-3, SM-AA-4, SM-CI-4, SM-ME-3). The discovery of pyrimidine-rich motifs was
 120 unexpected, since the treatments have been reported to primarily target purines [5, 24, 25, 26]. However, these findings

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

121 align with a previous report [23] that the ionic current is strongly influenced by the nucleotide in the seventh pore
122 position in the nanopore we used in our experiments (ONT R10.4.1). For signal variability, motifs of repeated purines
123 showed a strong effect in all treatments (Fig. 4A; SD-AA-1, SD-CI-2, SD-CI-3, SD-ME-1), while cisplatin also yielded
124 a thymine repeat (SD-CI-1), with similar patterns seen in AAI (SD-AA-2, SD-AA-4).

125 Strong motifs associating with dwell times were observed in AAI and melphalan (DT-AA-1, DT-ME-1; Fig. 4B). These
126 purine-rich motifs resembled those identified in signal SD analyses (SD-AA-1, SD-ME-1). Motifs linked to dwell
127 time differences at the motor protein (ADT) revealed an adenine-rich motif in AAI (MT-AA-1) and a less specific
128 motif in melphalan (MT-ME-1). Such motifs may indicate sequence contexts where adducts influence translocation
129 speed by interacting with the motor protein rather than the sensing region. In general, dwell time motifs contained
130 homopolymeric tracts, suggesting that repetitive sequences are particularly sensitive to adduct-induced perturbations.

131 AEAD analysis, which sensitively captures effects on longer sequence contexts (Methods), highlighted the role of
132 purine-rich motifs. For the AAI samples, to focus on the most variable set of 9-mers, we used the 9-mers within
133 the highest 1% of CvM test scores. In AAI and melphalan, the signal mean AEAD varied the most at A-rich motifs
134 (AM-AA-1, AM-ME-1). In the melphalan samples, we observed also thymine-rich motifs (AM-ME-2, AM-ME-3).
135 These were the shortest and least position-specific motifs discovered.

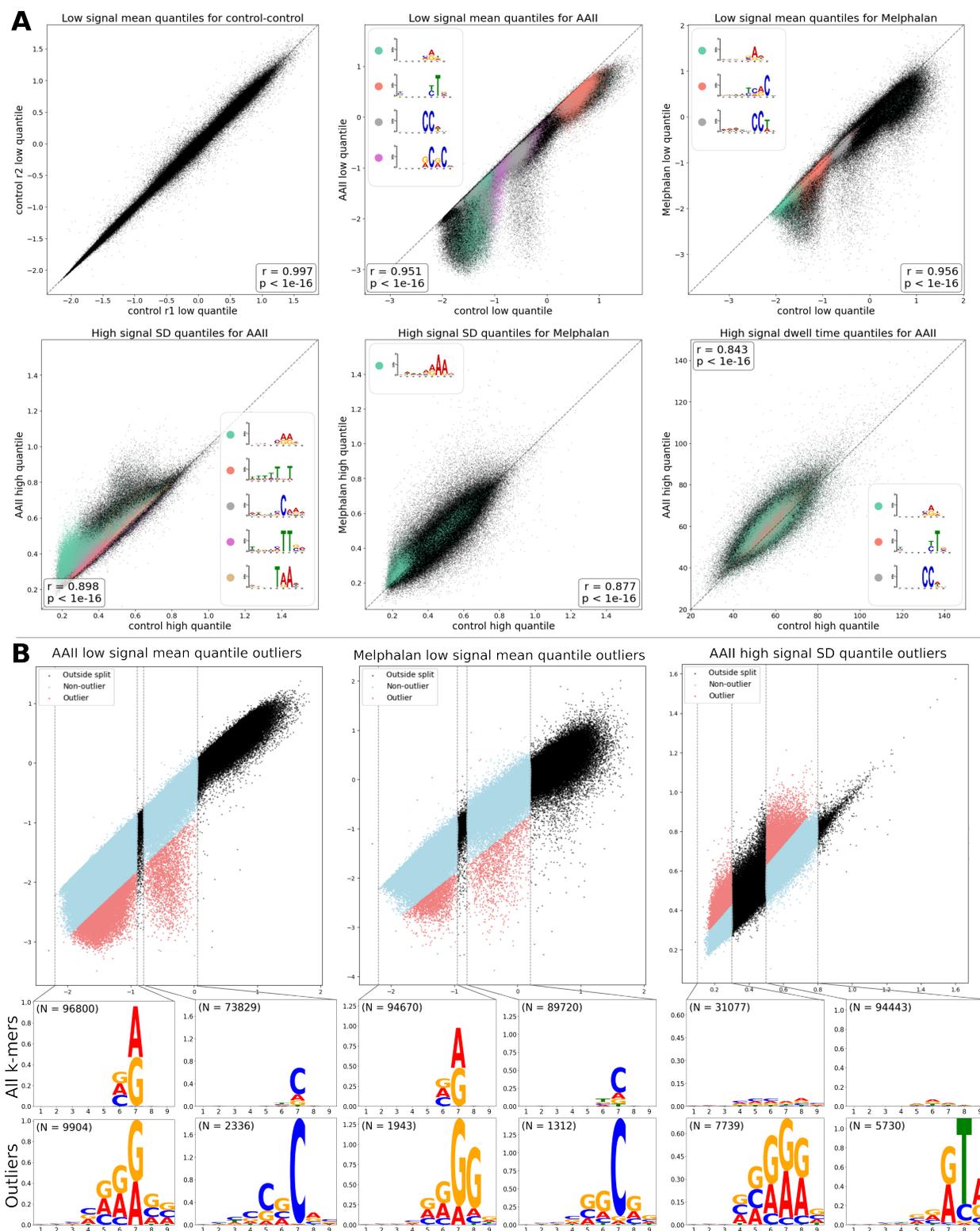
136 Finally, we found gapped sequence motifs associated with signal SD AEAD in every treatment group. In AAI, motif
137 AS-AA-1 is an A-rich motif, while in cisplatin and melphalan, the motifs contain adenines and thymines (AS-CI-1,
138 AS-ME-1). For mitomycin C, we found T-rich motifs (AS-CE-1, AS-CR-1), and for the interstrand mitomycin C
139 sample, an additional purine-rich motif, AS-CE-2. Signal SD AEAD was the only variable associated with significant
140 motifs across all treatments.

141 2.5 Recurrent sequence motifs display characteristic effects on nanopore signal levels

142 Quantile analyses revealed that adduct-induced effects were concentrated at the extremes of k -mer distributions, but
143 their directionality and recurrent sequence patterns remained unclear. We thus examined sequence motifs associated
144 with extreme quantiles in more detail (Methods).

145 Signal levels in both untreated and mitomycin-treated samples correlated closely with controls (control-control $r =$
146 0.997, MMC interstrand $r = 0.999$, and MMC intrastrand $r = 0.999$; Fig. 5A, Suppl. Fig. 5), indicating the absence
147 of strong effects in these conditions. In contrast, in both AAI and melphalan samples, nearly all k -mers displayed
148 reduced signal levels compared to untreated controls. Moreover, motifs discovered earlier were consistently associated
149 with lower signal levels in AAI, but the association was not as strong in melphalan samples (Fig. 5A). Examining the
150 k -mers which yielded the largest reductions revealed patterns (Fig. 5B) compatible with the previously discovered
151 signal mean and SD motifs (Fig. 4A). In melphalan samples, we observed that motifs associated with reduced signal
152 levels contained more guanines at the seventh pore position and more purines at the preceding positions, which was
153 also observed in AAI. Outliers at intermediate signal levels contained more cytosines at the seventh position and

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats



Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

154 showed increased frequencies of purines or cytosines in the preceding positions in both the AAI- and melphalan-treated
155 samples. Finally, for AAI, SD outliers at low signal values were enriched in purines at pore positions 6–8, and the high
156 signal SD outliers had a [A/G][T/C][A/G] motif at the end of the pore. In contrast, in AAI, the dwell time motifs did
157 not consistently associate with quantile differences, suggesting that this treatment has no consistent effect on nanopore
158 sequencing translocation speed.

159 We also tested how different DNA triplets affect extreme values of the sequencing statistics at specific nanopore
160 positions (Suppl. Fig. 6, 7). We found that purine-rich triplets were associated with lower signal levels in AAI and
161 melphalan samples. More specifically, a CAG-triplet in AAI and a GGG-triplet in melphalan had the highest effects.
162 In AAI samples, purine-rich triplets were associated with high signal variation, with AGG and GGG showing the most
163 variability.

164 In summary, motif discovery revealed sequence logos that were significantly associated with changes in the ionic
165 current. These motifs had purine bases, primarily adenines, at the upper pore positions, in all samples and for all signal
166 statistics, except for SD in cisplatin. In addition, we observed pyrimidine-rich motifs that should not be binding sites of
167 any of the adduct-forming compounds [5, 24, 25, 26]. The discovered motifs correlated with low signal mean quantiles
168 in AAI and melphalan treatments and with high standard deviation quantiles in AAI. Signal SD AEAD was the only
169 variable that had significant sequence motifs associated with every treatment, and the motifs had a unique gapped
170 appearance.

171 **2.6 Read interruptions are associated with upstream guanine repeats at the motor protein**

172 Nanopore sequencing occasionally yields incomplete reads due to interruptions in pore translocation [27]. By examining
173 rear barcode qualities (Methods), we observed that AAI, cisplatin, and melphalan samples produced more reads
174 lacking a rear barcode than control and MMC samples, indicating a higher frequency of interrupted reads (Fig. 6A,
175 Suppl. Fig. 1). To investigate whether these interruptions are associated with particular sequence contexts, we performed
176 motif discovery analysis on interrupted reads. We identified two types of interrupted reads based on the signal level at
177 the read end: passing-like reads where the level was similar to fully-read reads, and blocking reads where the level was
178 lower (Fig. 6B).

179 We discovered guanine-rich motifs in AAI and melphalan samples associating with read interruptions ($E < 0.05$; Fig.
180 6C). Passing-like reads in melphalan sample highlighted a simple guanine motif, whereas two blocking read motifs in
181 AAI and melphalan samples had more varied T(G) n and GGC compositions, respectively. Interestingly, these three
182 motifs were all enriched at the motor protein instead of the sensing region (Fig. 6D). This led us to believe that both
183 AAI and melphalan monoadducts may cause interruptions in nanopore sequencing. The GGC-motif in melphalan
184 blocking reads could suggest a binding position for an interstrand crosslink, which would interrupt sequencing by
185 interacting with the motor protein.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

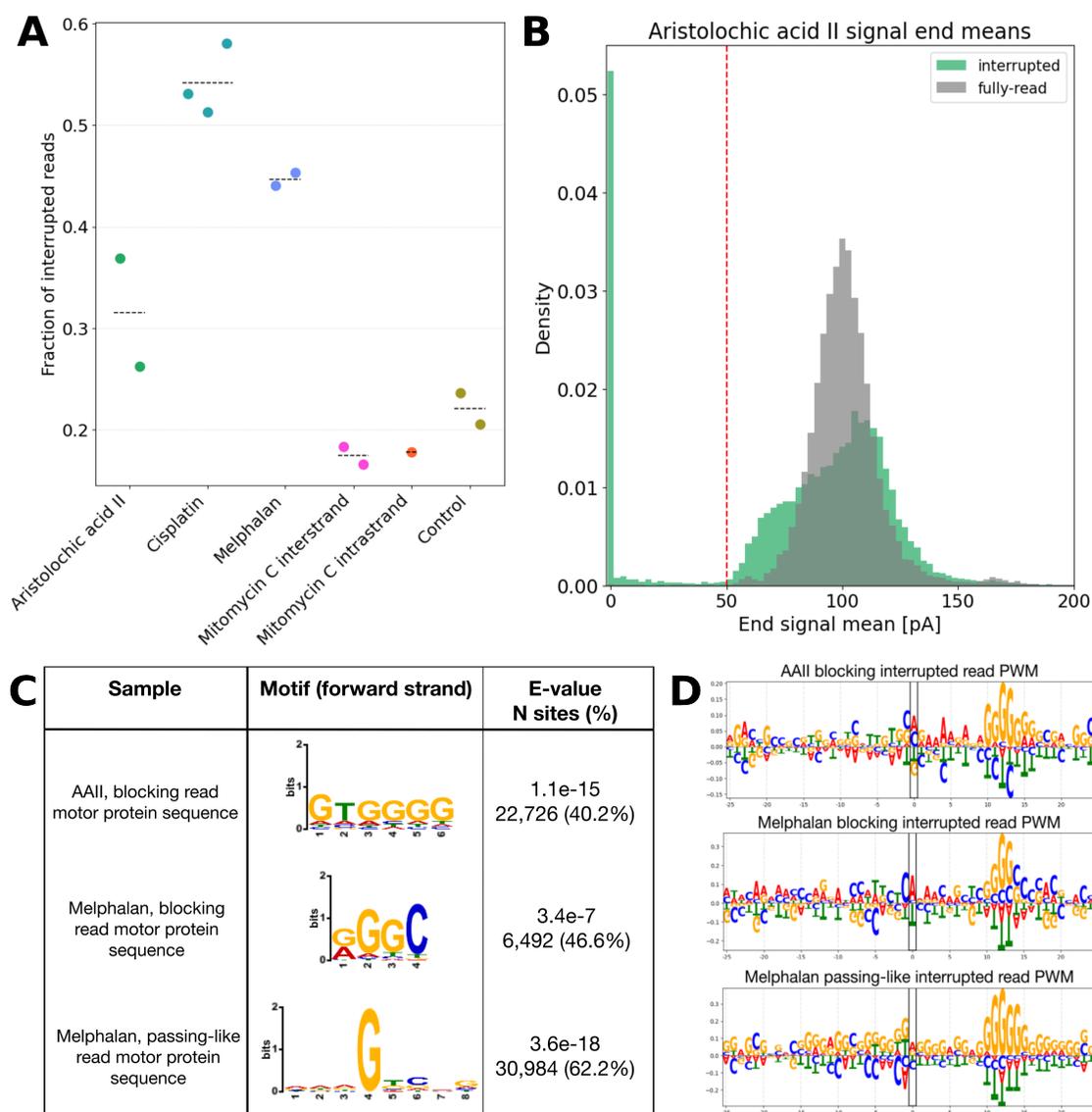


Figure 6: **A:** Fraction of interrupted reads in each sample replicate. **B:** Mean signal levels at read ends for interrupted and fully-read reads in the AAIL sample. The red line shows the classification threshold for passing-like and blocking reads. **C:** Sequence motifs detected with STREME associating with read interruptions. **D:** Position weight matrices (PWMs) for the sequences corresponding to motifs centered at the sensing region. The motor protein region is located approximately at position +14 from the sensing region.

186 3 Discussion

187 In this study, we investigated how genotoxic, adduct-forming compounds influence nanopore sequencing readouts. To
 188 our knowledge, this was the first study to thoroughly examine the effect of DNA adducts in complex DNA contexts
 189 (99.96% of all 9-mers), whereas previous studies have analyzed only up to four specific adduct sites [19, 20]. By
 190 applying our novel IonStats framework, we profiled adduct-induced changes across multiple dimensions of nanopore
 191 data, including base quality, nanopore signal statistics, pore translocation (dwell) times, and read interruptions. This
 192 effort revealed a multitude of compound-specific alterations in signal properties and sequence contexts, highlighting

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

193 both shared and distinct patterns of how different adducts perturb DNA translocation through the pore. Our observations
194 establish a foundation for linking nanopore-based adduct detection to exposure-associated mutational processes and
195 signatures [28, 3].

196 Detecting DNA adducts is a crucial step toward understanding how mutational signatures arise. For instance, signature
197 SBS4 found in lung cancers has been linked to tobacco carcinogens [29]. Site-directed mutagenesis studies demonstrate
198 that the mutational signatures of DNA adducts are directly related to DNA polymerase activity, its DNA insertion and
199 extension mechanism, and the DNA sequence context [30]. To properly assign an etiology to a mutational signature,
200 it is essential to consider DNA adduct types and locations as the key intermediate steps that yield somatic mutations.
201 This, however, requires accurate, site-specific methods for adduct detection: we here showed that nanopore sequencing
202 coupled with sophisticated computational tools for ionic current analysis shows substantial promise.

203 Our software package IonStats can be easily applied to sets of treated and control samples to evaluate sequencing
204 statistics for treatment effects. In contrast to previous DNA and RNA modification detection tools that can be used to
205 characterize DNA adducts [31, 19, 32], IonStats was developed specifically for DNA adduct detection. It introduces
206 several new readouts on signal-level statistics, and translocation dynamics including dwell-times and read interruptions.
207 These readouts enable multifaceted characterization of effects caused by DNA adducts. We hope that this tool
208 will encourage further research on profiling effects of environmental exposures with nanopore sequencing beyond
209 metagenomics [33].

210 Our treatment-specific findings illustrate both potential and challenges. Aristolochic acid II (AAII) altered sequencing
211 quality and ionic current signals mainly in purine-rich contexts at the nanopore sensing region. In addition, guanine
212 repeats were enriched upstream of the sensing region at the motor protein, associating with translocation interruptions.
213 This finding led us to believe that aristolactam II-dA (AL-dA) adducts are more likely to pass through the nanopore,
214 while AL-dG adducts can cause read interruptions by interacting with the motor protein. This is consistent with the
215 larger size of AL-dG adducts compared to AL-dA adducts [5]. Melphalan produced shorter and poorer quality reads
216 than the other treatments, and displayed widespread effects in adenine-rich contexts and evidence of read interruptions
217 in CG-rich contexts. This could indicate melphalan crosslinks, binding to two guanines on opposing strands [34], to
218 cause translocation failure at the motor protein, similarly to AAII.

219 We observed cisplatin-driven differences in signal statistics primarily in A and T-rich motifs. Although GG-crosslinks
220 are the primary reported cisplatin adduct, our motif discovery analysis was not able to detect these. This may be due
221 to insufficient power resulting from the low sequencing coverage, caused by degradation of the DNA by the cisplatin
222 treatment. We were surprised to find many T-rich motifs in cisplatin and AAII (*e.g.*, SM-AA-2, SM-CI-2, SD-AA-2,
223 SD-AA-4, SD-CI-1), since cisplatin has not been reported to bind to thymines, although GpNpG 1,3-intrastrand
224 crosslink can rarely occur (SM-CI-3) [35]. Furthermore, AL-dT adducts have not been observed in previous studies [36].
225 The only effects we discovered after MMC treatments were the motifs associated with signal variability over longer
226 sequences captured by the AEAD score, suggesting that AEAD is more sensitive to treatment effects than other statistics
227 studied.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

228 While IonStats enabled comprehensive profiling of treatment effects, it is not able to identify individual reads harboring
229 adducts. Thus, its sensitivity might be limited when adduct frequencies are low. In addition, IonStats has only been
230 applied to *in vitro* samples that are expected to carry relatively large amounts of DNA adducts compared to *in vivo*
231 samples [3]. Further work is required to evaluate the utility of IonStats in *in vivo* settings. Another limitation of our
232 work is that IonStats is not able to quantify the amount of adducts in treated samples directly. A complementary mass
233 spectrometry analysis could be used to estimate the fraction of adducted nucleotides.

234 Taken together, our work establishes a framework for studying DNA adducts with nanopore sequencing and demonstrates
235 how different compounds leave distinct signatures in sequencing signals. By developing methods for site-specific
236 adduct detection, future studies may bridge the gap between DNA lesions and mutational signatures, clarifying the
237 causal pathways by which exogenous and endogenous exposures contribute to cancer development.

238 References

- 239 [1] Byeong Hwa Yun, Jingshu Guo, Medjda Bellamri, and Robert J. Turesky. DNA adducts: Formation, biological
240 effects, and new biospecimens for mass spectrometric measurements in humans. *Mass Spectrometry Reviews*,
241 39(1-2):55–82, 2020.
- 242 [2] Aziz Sancar, Laura A. Lindsey-Boltz, Keziban Ünsal-Kaçmaz, and Stuart Linn. Molecular mechanisms of
243 mammalian DNA repair and the DNA damage checkpoints. *Annual Review of Biochemistry*, 73:39–85, 2004.
- 244 [3] Gunnar Boysen, Ludmil B. Alexandrov, Raheleh Rahbari, Intawat Nookaew, Dave Ussery, Mu Rong Chao,
245 Chiung Wen Hu, and Marcus S. Cooke. Investigating the origins of the mutational signatures in cancer. *Nucleic
246 Acids Research*, 53(1), 2025.
- 247 [4] Thomas A. Rosenquist and Arthur P. Grollman. Mutational signature of aristolochic acid: Clue to the recognition
248 of a global disease. *DNA Repair*, 44:205–211, 2016.
- 249 [5] Marie Stiborová, Volker M. Arlt, and Heinz H. Schmeiser. DNA adducts formed by aristolochic acid are unique
250 biomarkers of exposure and explain the initiation phase of upper urothelial cancer. *International Journal of
251 Molecular Sciences*, 18(10), 2017.
- 252 [6] Michael Alavanja, John A. Baron, Ross C. Brownson, Patricia A. Buffler, David M. DeMarini, Mirjana V.
253 Djordjevic, Richard Doll, Elizabeth T.H. Fontham, Yu Tang Gao, Nigel Gray, Prakash C. Gupta, Allan Hackshaw,
254 Stephen S. Hecht, Kirsti Husgafvel-Pursiainen, Elena Matos, Richard Peto, David H. Phillips, Jonathan M. Samet,
255 Gary Stoner, Michael J. Thun, Jean Trédaniel, Paolo Vineis, H. Erich Wichmann, Anna H. Wu, and David Zaridze.
256 Tobacco smoke and involuntary smoking. *IARC Monographs on the Evaluation of Carcinogenic Risks to Humans*,
257 83:1–1413, 2004.
- 258 [7] Thomas W. Kensler, Bill D. Roebuck, Gerald N. Wogan, and John D. Groopman. Aflatoxin: A 50-year Odyssey
259 of mechanistic and translational toxicology. *Toxicological Sciences*, 120(SUPPL.1):28–48, 2011.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

- 260 [8] M. Austin Argentieri, Najaf Amin, Alejo J. Nevado-Holgado, William Sproviero, Jennifer A. Collister, Sarai M.
261 Keestra, Midas M. Kuilman, Bigina N.R. Ginos, Mohsen Ghanbari, Aiden Doherty, David J. Hunter, Alexandra
262 Alvergne, and Cornelia M. van Duijn. Integrating the environmental and genetic architectures of aging and
263 mortality. *Nature Medicine*, 31(3):1016–1025, 2025.
- 264 [9] Mengqiu Cao and Xinyu Zhang. DNA Adductomics: A Narrative Review of Its Development, Applications, and
265 Future. *Biomolecules*, 14(9), 2024.
- 266 [10] David K. La and James A. Swenberg. DNA adducts: Biological markers of exposure and potential applications to
267 risk assessment. *Mutation Research - Reviews in Genetic Toxicology*, 365(1-3):129–146, 1996.
- 268 [11] Daniel R. Semlow and Johannes C. Walter. Mechanisms of Vertebrate DNA Interstrand Cross-Link Repair.
269 *Annual Review of Biochemistry*, 90:107–135, 2021.
- 270 [12] Natalia Y. Tretyakova, Arnold Groehler, and Shaofei Ji. DNA-Protein Cross-Links: Formation, Structural
271 Identities, and Biological Outcomes. *Accounts of Chemical Research*, 48(6):1631–1644, 2015.
- 272 [13] Romel P. Dator, Kevin J. Murray, Matthew W. Luedtke, Foster C. Jacobs, Fekadu Kassie, Hai Dang Nguyen,
273 Peter W. Villalta, and Silvia Balbo. Identification of Formaldehyde-Induced DNA-RNA Cross-Links in the A/J
274 Mouse Lung Tumorigenesis Model. *Chemical Research in Toxicology*, 35(11):2025–2036, 2022.
- 275 [14] Gunnar Boysen and Intawat Nookaew. Current and Future Methodology for Quantitation and Site-Specific
276 Mapping the Location of DNA Adducts. *Toxics*, 10(2), feb 2022.
- 277 [15] Lieselot Y. Hemeryck, Anneleen I. Decloedt, Julie Vanden Bussche, Karen P. Geboes, and Lynn Vanhaecke. High
278 resolution mass spectrometry based profiling of diet-related deoxyribonucleic acid adducts. *Analytica Chimica*
279 *Acta*, 892:123–131, 2015.
- 280 [16] Yuan Jhe Chang, Marcus S. Cooke, Chiung Wen Hu, and Mu Rong Chao. Novel approach to integrated
281 DNA adductomics for the assessment of in vitro and in vivo environmental exposures. *Archives of Toxicology*,
282 92(8):2665–2680, 2018.
- 283 [17] Mikhail F Denissenko, Sundaresan Venkatachalam, Yu Hua Ma, and Altaf A Wani. Site-specific induction and
284 repair of benzo[a]pyrene diol epoxide DNA damage in human H-ras protooncogene as revealed by restriction
285 cleavage inhibition. *Mutation Research - DNA Repair*, 363(1):27–42, 1996.
- 286 [18] Bo Cao, Xiaolin Wu, Jieliang Zhou, Hang Wu, Lili Liu, Qinghua Zhang, Michael S Demott, Chen Gu, Lian-
287 rong Wang, Delin You, and Peter C Dedon. Nick-seq for single-nucleotide resolution genomic maps of DNA
288 modifications and damage. *Nucleic Acids Research*, 48(12):6715–6725, 2020.
- 289 [19] Intawat Nookaew, Gunnar Boysen, Piroon Jenjaroenpun, Hua Du, Pengcheng Wang, Jun Wu, Thidathip Wong-
290 surawat, Sun Hee Moon, En Huang, and Yinsheng Wang. Detection and discrimination of DNA adducts
291 differing in size, regiochemistry, and functional group by nanopore sequencing. *Chemical Research in Toxicology*,
292 33(12):2944–2952, dec 2020.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

- 293 [20] Xinjia Zhao, Yuru Liu, Xiaoyu Chen, Zhuang Mi, Wei Li, Pengye Wang, Xinyan Shan, and Xinghua Lu.
294 Detection and Characterization of Single Cisplatin Adducts on DNA by Nanopore Sequencing. *ACS Omega*,
295 6(26):17027–17034, jul 2021.
- 296 [21] Yunhao Wang, Yue Zhao, Audrey Bollas, Yuru Wang, and Kin Fai Au. Nanopore sequencing technology,
297 bioinformatics and applications. *Nature Biotechnology*, 39(11):1348–1365, nov 2021.
- 298 [22] Miten Jain, Hugh E. Olsen, Benedict Paten, and Mark Akeson. The Oxford Nanopore MinION: delivery of
299 nanopore sequencing to the genomics community. *Genome Biology*, 17(1), dec 2016.
- 300 [23] Sam Kovaka, Paul W. Hook, Katharine M. Jenike, Vikram Shivakumar, Luke B. Morina, Roham Razaghi, Winston
301 Timp, and Michael C. Schatz. Uncalled4 improves nanopore DNA and RNA modification detection via fast and
302 accurate signal alignment. *Nature Methods*, 22(4):681–691, 2025.
- 303 [24] Shaloam Dasari and Paul Bernard. Cisplatin in cancer therapy : Molecular mechanisms of action. *European*
304 *Journal of Pharmacology*, 740:364–378, 2014.
- 305 [25] M R Osborne and P. D. Lawley. Alkylation of DNA by melphalan with special reference to adenine derivatives
306 and adenine-guanine cross-linking. *Chemico-Biological Interactions*, 89(1):49–60, 1993.
- 307 [26] Maria Tomasz and Yolanda Palom. The mitomycin bioreductive antitumor agents: Cross-linking and alkylation of
308 DNA as the molecular basis of their activity, 1997.
- 309 [27] D. Branton and D.W. Deamer. *Nanopore Sequencing: An Introduction*. World Scientific Publishing Company,
310 2019.
- 311 [28] Lauri J. Sipilä, Riku Katainen, Mervi Aavikko, Janne Ravanti, Iikki Donner, Rainer Lehtonen, Ilmo Leivo, Henrik
312 Wolff, Reetta Holmila, Kirsti Husgafvel-Pursiainen, and Lauri A. Aaltonen. Genome-wide somatic mutation
313 analysis of sinonasal adenocarcinoma with and without wood dust exposure. *Genes and Environment*, 46(1):1–13,
314 2024.
- 315 [29] Ludmil B. Alexandrov, Serena Nik-Zainal, David C. Wedge, Samuel A.J.R. Aparicio, Sam Behjati, Andrew V.
316 Biankin, Graham R. Bignell, Niccolò Bolli, Ake Borg, Anne Lise Børresen-Dale, Sandrine Boyault, Birgit
317 Burkhardt, Adam P. Butler, Carlos Caldas, Helen R. Davies, Christine Desmedt, Roland Eils, Jórunn Erla Eyfjörd,
318 John A. Foekens, Mel Greaves, Fumie Hosoda, Barbara Hutter, Tomislav Ilicic, Sandrine Imbeaud, Marcin
319 Imielinski, Natalie Jäger, David T.W. Jones, David Jonas, Stian Knappskog, Marcel Koo, Sunil R. Lakhani,
320 Carlos López-Otín, Sancha Martin, Nikhil C. Munshi, Hiromi Nakamura, Paul A. Northcott, Marina Pajic, Elli
321 Papaemmanuil, Angelo Paradiso, John V. Pearson, Xose S. Puente, Keiran Raine, Manasa Ramakrishna, Andrea L.
322 Richardson, Julia Richter, Philip Rosenstiel, Matthias Schlesner, Ton N. Schumacher, Paul N. Span, Jon W.
323 Teague, Yasushi Totoki, Andrew N.J. Tutt, Rafael Valdés-Mas, Marit M. Van Buuren, Laura Van 'T Veer, Anne
324 Vincent-Salomon, Nicola Waddell, Lucy R. Yates, Jessica Zucman-Rossi, P. Andrew Futreal, Ultan McDermott,
325 Peter Lichter, Matthew Meyerson, Sean M. Grimmond, Reiner Siebert, Elías Campo, Tatsuhiro Shibata, Stefan M.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

- 326 Pfister, Peter J. Campbell, and Michael R. Stratton. Signatures of mutational processes in human cancer. *Nature*,
327 500(7463):415–421, 2013.
- 328 [30] Robert P. Fuchs and Shingo Fujii. Translesion DNA synthesis and mutagenesis in prokaryotes. *Cold Spring*
329 *Harbor Perspectives in Biology*, 5(12):1–22, 2013.
- 330 [31] Piroon Jenjaroenpun, Thidathip Wongsurawat, Taylor D Wadley, Trudy M Wassenaar, Jun Liu, Qing Dai, Visanu
331 Wanchai, Nisreen S Akel, Azemat Jamshidi-Parsian, Aime T Franco, Gunnar Boysen, Michael L Jennings,
332 David W Ussery, Chuan He, and Intawat Nookaew. Decoding the epitranscriptional landscape from native RNA
333 sequences. *Nucleic Acids Research*, 49(2):1–13, 2021.
- 334 [32] Jared T. Simpson, Rachael E. Workman, P. C. Zuzarte, Matei David, L. J. Dursi, and Winston Timp. Detecting
335 DNA cytosine methylation using nanopore sequencing. *Nature Methods*, 14(4):407–410, 2017.
- 336 [33] Mantas Sereika, Rasmus Hansen Kirkegaard, Søren Michael Karst, Thomas Yssing Michaelsen, Emil Aarre
337 Sørensen, Rasmus Dam Wollenberg, and Mads Albertsen. Oxford Nanopore R10.4 long-read sequencing enables
338 the generation of near-finished bacterial genomes from pure cultures and metagenomes without short-read or
339 reference polishing. *Nature Methods*, 19(7):823–826, jul 2022.
- 340 [34] Lawrence F Povirk and David E Shuker. DNA damage and mutagenesis induced by nitrogen mustards, 1994.
- 341 [35] Askar Yimit, Ogun Adebali, Aziz Sancar, and Yuchao Jiang. Differential damage and repair of DNA-adducts
342 induced by anti-cancer drug cisplatin across mouse organs. *Nature Communications*, 10(1), 2019.
- 343 [36] Wan Chan, Yufang Zheng, and Zongwei Cai. Liquid Chromatography-Tandem Mass Spectrometry Analysis of
344 the DNA Adducts of Aristolochic Acids. *Journal of the American Society for Mass Spectrometry*, 18(4):642–650,
345 2007.
- 346 [37] Volker M Arlt, Heinz H Schmeiser, and Gerd P Pfeifer. Sequence-specific detection of aristolochic acid-DNA
347 adducts in the human p53 gene by terminal transferase-dependent PCR The plant extract aristolochic acid (AA), a
348 mixture consisting. Technical Report 1, 2001.
- 349 [38] Fubo Ma, Shuanghong Yan, Jinyue Zhang, Yu Wang, Liying Wang, Yuqin Wang, Shanyu Zhang, Xiaoyu Du,
350 Panke Zhang, Hong-Yuan Chen, and Shuo Huang. Nanopore Sequencing Accurately Identifies the Cisplatin
351 Adduct on DNA. *ACS Sensors*, 2021.
- 352 [39] Dalia Mohamed and Michael Linscheid. Separation and identification of trinucleotide-melphalan adducts from
353 enzymatically digested DNA using HPLC-ESI-MS. In *Analytical and Bioanalytical Chemistry*, volume 392,
354 pages 805–817, nov 2008.
- 355 [40] Manuel M. Paz, Sweta Ladwa, Elise Champeil, Yanfeng Liu, Sara Rockwell, Ernest K. Boamah, Jill Bargonetti,
356 John Callahan, John Roach, and Maria Tomasz. Mapping DNA adducts of mitomycin C and decarbamoyl
357 mitomycin C in cell lines using liquid chromatography/electrospray tandem mass spectrometry. *Chemical*
358 *Research in Toxicology*, 21(12):2370–2378, dec 2008.
- 359 [41] Heng Li. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18):3094–3100, 2018.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

- 360 [42] Jurgen F. Nijkamp, Marcel van den Broek, Erwin Datema, Stefan de Kok, Lizanne Bosman, Marijke A. Luttk, 361 Pascale Daran-Lapujade, Wanwipa Vongsangnak, Jens Nielsen, Wilbert H.M. Heijne, Paul Klaassen, Chris J. 362 Paddon, Darren Platt, Peter Kötter, Roeland C. van Ham, Marcel J.T. Reinders, Jack T. Pronk, Dick de Ridder, 363 and Jean Marc Daran. De novo sequencing, assembly and analysis of the genome of the laboratory strain 364 *Saccharomyces cerevisiae* CEN.PK113-7D, a model for modern industrial biotechnology. *Microbial Cell 365 Factories*, 11(1):36, 2012.
- 366 [43] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni 367 Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua 368 Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R.J. Nelson, Eric Jones, Robert Kern, Eric Larson, 369 C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert 370 Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian 371 Pedregosa, Paul van Mulbregt, Aditya Vijaykumar, Alessandro Pietro Bardelli, Alex Rothberg, Andreas Hilboll, 372 Andreas Kloeckner, Anthony Scopatz, Antony Lee, Ariel Rokem, C. Nathan Woods, Chad Fulton, Charles 373 Masson, Christian Häggström, Clark Fitzgerald, David A. Nicholson, David R. Hagen, Dmitrii V. Pasechnik, 374 Emanuele Olivetti, Eric Martin, Eric Wieser, Fabrice Silva, Felix Lenders, Florian Wilhelm, G. Young, Gavin A. 375 Price, Gert Ludwig Ingold, Gregory E. Allen, Gregory R. Lee, Hervé Audren, Irvin Probst, Jörg P. Dietrich, 376 Jacob Silterra, James T. Webber, Janko Slavič, Joel Nothman, Johannes Buchner, Johannes Kulick, Johannes L. 377 Schönberger, José Vinícius de Miranda Cardoso, Joscha Reimer, Joseph Harrington, Juan Luis Cano Rodríguez, 378 Juan Nunez-Iglesias, Justin Kuczynski, Kevin Tritz, Martin Thoma, Matthew Newville, Matthias Kümmerer, 379 Maximilian Bolingbroke, Michael Tartre, Mikhail Pak, Nathaniel J. Smith, Nikolai Nowaczyk, Nikolay Shebanov, 380 Oleksandr Pavlyk, Per A. Brodtkorb, Perry Lee, Robert T. McGibbon, Roman Feldbauer, Sam Lewis, Sam Tygier, 381 Scott Sievert, Sebastiano Vigna, Stefan Peterson, Surhud More, Tadeusz Pudlik, Takuya Oshima, Thomas J. 382 Pingel, Thomas P. Robitaille, Thomas Spura, Thouis R. Jones, Tim Cera, Tim Leslie, Tiziano Zito, Tom Krauss, 383 Utkarsh Upadhyay, Yaroslav O. Halchenko, and Yoshiki Vázquez-Baeza. SciPy 1.0: fundamental algorithms for 384 scientific computing in Python. *Nature Methods*, 17(3):261–272, 2020.
- 385 [44] Timothy L. Bailey. STREME: accurate and versatile sequence motif discovery. *Bioinformatics*, 37(18):2834–2840, 386 2021.
- 387 [45] Ammar Tareen and Justin B Kinney. Logomaker: Beautiful sequence logos in Python. *Bioinformatics*, 36(7):2272– 388 2274, 2020.
- 389 [46] Fabian A. Buske, Mikael Bodén, Denis C. Bauer, and Timothy L. Bailey. Assigning roles to DNA regulatory 390 motifs using comparative genomics. *Bioinformatics*, 26(7):860–866, 2010.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

391 **Acknowledgments**

392 We would like to thank Veera Erkkilä, Pinja Perkkiö, Iina Vuoristo, and Inga-Lill Åberg for technical support, and Harri
393 Kangas for discussions. We acknowledge CSC – IT Center for Science, Finland, for generous computational resources.
394 Figure 1D was created with BioRender.com. This study was supported by the Research Council of Finland (#328890 to
395 EP) and the University of Helsinki Doctoral Programme in Integrative Life Science.

396 **Data availability statement**

397 All sequencing data for this study have been deposited in the European Nucleotide Archive (ENA) at EMBL-EBI under
398 accession number PRJEB101972.

399 **4 Methods**

400 **4.1 Whole genome DNA amplification**

401 Yeast genomic DNA was extracted from the haploid *Saccharomyces cerevisiae* strain CEN.PK113-7D by using the
402 MasterPure Yeast DNA Purification Kit (MPY80200, LGC). High molecular weight yeast genomic DNA was uniformly
403 amplified by using the REPLI-G mini whole genome amplification kit (150023, Qiagen). In brief, 20 ng of yeast
404 genomic DNA was used as a template and amplified by incubating with REPLI-G DNA Polymerase at 30 °C for 16
405 hours. Amplified yeast genomic DNA was purified by using 1.8x Ampure XP beads (10136224, Beckman Coulter)
406 and eluted in 10 mM Tris-HCl, pH 8.0, or double-distilled water. Amplified yeast genomic DNA was fragmented
407 and size-selected for 10-15 Kb by using g-tube (520079, Covaris). For Melphalan treatment, DNA was extracted for
408 ultrapure distilled water (10977035, Invitrogen).

409 **4.2 Aristolochic acid II treatment**

410 AAI DNA adducts were formed by using the protocol described in [37]. Amplified yeast genomic DNA (10 µg DNA,
411 one µg per µl) was treated with 0.25 mM AAI solution in 50 µl of 50 mM potassium phosphate buffer, pH 5.8. AAI
412 was activated by 1 mg of zinc dust and incubated at 37 °C for 30 minutes. AAI-treated DNA was purified by using
413 1.8x Ampure XP beads and eluted in 25 µl 10 mM Tris-HCl, pH 8.0.

414 **4.3 Cisplatin treatment**

415 Cisplatin DNA-adduct was prepared by using the protocol in [38] with modifications. Cisplatin (3 mg, 10 µmol) was
416 activated by incubation with AgNO₃ (6 mg, 35 µmol) at 37 °C for 6 hours in 50 µl of water. The activation reaction
417 was conducted in a 1.5 ml Eppendorf tube covered with an aluminium foil to prevent light exposure. Cisplatin and
418 AgNO₃ mixture was centrifuged at 13,000 rpm for 10 minutes, and the clear supernatant containing activated cisplatin
419 was used for DNA-adduct formation. Cisplatin-DNA adduct was formed by incubating DNA (10 µg) and activated

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

420 cisplatin in 10 mM HEPES buffer, pH 7.0, at 37 °C for 1 hour in the dark. Cisplatin-treated DNA was purified by using
421 Phenol-Chloroform extraction and eluted in 10 mM Tris, pH 8.0.

422 **4.4 Melphalan treatment**

423 Melphalan DNA-adduct was prepared by using the protocol in [39] with a modification. Amplified yeast genomic
424 DNA was eluted in distilled water (10 µg in 10 µl water). Amplified yeast genomic DNA was incubated with 8 µl of
425 melphalan stock solution (1 mg per ml of DMSO) and 12 µl of DMSO at 37 °C for 1 hour. Melphalan-treated DNA was
426 purified by using Phenol-Chloroform extraction and eluted in distilled water.

427 **4.5 Mitomycin C treatment**

428 **4.5.1 Intrastrand crosslink adduct preparation:**

429 Mitomycin C (MMC) intrastrand crosslink DNA-adduct was prepared by using the protocol in [40]. Amplified yeast
430 genomic DNA was diluted to 30 µl of 0.1 M PBS pH 7.4 (P5493, Sigma) (6.54 µg DNA in 30 µl PBS, 60 nmol of DNA).
431 MMC was dissolved in 0.1 M PBS, pH 7.4, to prepare a 100 mM stock solution. 0.24 µmol of mitomycin C (2.4 µl of
432 0.1 M MMC stock solution was diluted to 27 µl in 0.1 M PBS) and 12 µmol of Na₂S₂O₄ (3 µl of 40 mM Na₂S₂O₄ in
433 water) were added to amplified yeast genomic DNA and incubated in a PCR tube in a thermoshaker at 40 °C, 400 rpm
434 with open lid in fume hood. MMC-treated DNA was purified by using 1.8x Ampure XP beads and eluted in 20 µl 10
435 mM Tris-HCl, pH 8.0.

436 **4.5.2 Interstrand crosslink adduct preparation:**

437 MMC interstrand crosslink DNA adduct was prepared using a protocol similar to intrastrand crosslinking, except the
438 incubation was performed in a closed lid PCR tube, and a total of 60 µmol of Na₂S₂O₄ was added in five portions of 12
439 µmol at a five-minute interval. MMC-treated DNA was purified by using 1.8x Ampure XP beads and eluted in 20 µl 10
440 mM Tris-HCl, pH 8.0.

441 **4.6 DNA sequencing and data collection**

442 Two replicates for each treatment were prepared for sequencing, with an additional cisplatin sample, where the DNA
443 fragments were selected by size. Approximately 400 ng of DNA was used for each sample to prepare the sequencing
444 library with the Native Barcoding Kit 24 V14 (SQK-NBD114.24, Oxford Nanopore Technologies (ONT)) following the
445 manufacturer's protocol. In total, 13 barcodes were used (two for each treatment, two for controls, and one additional
446 for cisplatin size-selected). The constructed libraries were loaded onto the R10.4.1/FLO-PRO114M flow cell on the
447 PromethION 24 device. The sequencing and data collection were performed with the MinKNOW software v24.02.19,
448 which generated .pod5 files containing the raw ionic signal measurements for individual DNA molecules.

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

449 The cisplatin treatment heavily damaged the DNA, and we observed a small number of reads for both non-size-selected
450 replicates and the size-selected replicate. For unknown reasons, the yield for the second MMC intrastrand replicate was
451 very low, and we excluded this sample from further analyses.

452 4.7 Data processing

453 Basecalling and alignment were performed with ONT's basecaller Dorado (version 0.6.0) with the
454 dna_r10.4.1_e8.2_400bps_sup@v4.3.0 model. The aligner function in Dorado uses minimap2 [41] to align the
455 basecalled reads to a reference genome with the "lr:hq" preset. We used the *S. cerevisiae* strain CEN.PK113-7D
456 assembly from [42] as the reference genome. To map the raw ionic current measurements to positions in the reference
457 genome, a signal-alignment step is required. The ionic current is divided into events, and each event is mapped to a
458 reference position. We applied Uncalled4 [23] to align signals to reference positions for all reads. This allowed us to
459 compute and compare ionic current statistics for different samples.

460 4.8 Data analysis of sequencing statistics with IonStats

461 IonStats can extract and compare quality scores, ionic current means, ionic current standard deviations, dwell times,
462 and adjusted dwell times between treated and control samples. We also derived windowed error scores for ionic current
463 mean and standard deviation, referred to here as Average Expected Absolute Deviation (AEAD) scores.

464 The quality score (Q) is a value given to each base by the basecaller, and it is logarithmically related to the estimated
465 basecalling error probability (P): $Q = -10 \log_{10} P$. The ionic current mean and standard deviation are computed
466 from the raw ionic current measurements that belong to a single event, defined by the signal alignment. The dwell time
467 refers to the time that a k -mer resides inside the nanopore, and we quantified it as the number of measurements that
468 belong to an event. To capture the interaction between putative adducts and nanopore motor protein, we computed the
469 motor-protein shifted dwell times (adjusted dwell times), where the dwell time of a k -mer was associated with a k -mer
470 14 base pairs downstream in the reference genome [20].

471 To estimate how much the ionic current varies over consecutive events, we propose the AEAD statistic. For a window
472 of width w , the AEAD for statistic s at position i can be specified as:

$$AEAD_{s,i} = \frac{1}{w} \sum_{j=i-\lfloor \frac{w}{2} \rfloor}^{i+\lfloor \frac{w}{2} \rfloor} |s_j - E[s, k_j]|$$

473 where s_j is the value of the statistic for event j and $E[s, k_j]$ is the expected value of the statistic s for k -mer centered at
474 position j . We used the 9-mer table provided by Uncalled4 to get the expected values for the signal mean and standard
475 deviation for all 9-mers.

476 IonStats collects values for selected variables and groups them by the reference k -mer or by reference position for
477 a specific chromosome or contig. It then performs distributional tests on the k -mer/ref_pos distributions to find the

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

478 k -mers that show statistically significant differences in the variables of interest. It applies the Cramér–von Mises (CvM)
479 criterion and the Kolmogorov–Smirnov (KS) test to compare the k -mers. For both of these tests, we used the versions
480 that were implemented in the SciPy [43] Python library. These two tests are applied to capture different effects since
481 KS is more sensitive to differences in medians, while CvM is more sensitive to differences in distribution tails. Since
482 both tests assume the data to be continuous, we added a small amount of noise to the discrete variables (quality score,
483 dwell time, and adjusted dwell time) to facilitate test calibration. To control for multiple testing, IonStats applies the
484 Benjamini-Hochberg correction for the p -values of each test with a significance level α . We used $\alpha = 0.05$.

485 **4.9 Motif discovery analysis**

486 To search for patterns that could correspond to recurrent adduct-forming sequence motifs in the DNA, we performed
487 motif discovery analysis with STREME [44] for the significant k -mers associated with each test. In the STREME
488 analyses, we used the set of all k -mers as the control sequences. To process the k -mers only in forward-strand
489 orientation, we used a custom STREME alphabet, which did not consider the complementarity of nucleobases. The
490 motif discovery was performed separately on forward and reverse reads. A motif had to occur on both strands with an
491 E-value smaller than 0.05 to be reported in Fig. 3 and considered significant. All motifs discovered by STREME with
492 $E < 0.05$ are shown in Suppl. Fig. 8-16. We also used the Logomaker [45] Python library to draw sequence logos for
493 the significant k -mers, and for the outlier k -mers in Fig. 5.

494 For the AAI signal mean AEAD, we restricted the motif discovery analysis only to the k -mers within the top 1% CvM
495 values. This was done to reduce the similarity between the positive set and the control set, which consisted of all k -mers.
496 Initially, by using all significant k -mers, we found no significant motifs, likely due to the high overlap of the positive
497 and control k -mer-sets in the STREME analysis.

498 **4.10 Quantile difference analysis**

499 To study how the extreme values of k -mer distributions compare in treated and control samples, we computed the
500 top and bottom 1%-quantiles for each sequencing variable. We then computed the quantile difference to measure if
501 a k -mer shows more extreme values in the treated samples. The quantile difference for k -mer i is defined simply as:
502 $d_i = q_{t,i} - q_{c,i}$, where c is control and t is treatment, and q can be either the lowest or highest 1%.

503 **4.11 Data analysis for interrupted reads**

504 Nanopore sequencing reads can exhibit low rear barcode qualities, indicating that the read barcode was non-determinable
505 based on the sequence at the end of the read. This led us to believe that interrupted reads, i.e., reads that interrupted
506 the sequencing and did not pass completely through the pore, belong to this subset of reads. We classified reads as
507 interrupted based on the `barcode_rear_score` value, which is recorded for all reads in the sequencing summary file
508 provided by MinKNOW and which ranges between -100 and 100 . If `barcode_rear_score` was < 60 , we classified

Compound-specific DNA adduct profiling with nanopore sequencing and IonStats

509 the read as interrupted. We chose this threshold as MinKNOW requires a minimum value of 60 from either the front or
510 rear barcode scores for a read to be assigned to a barcode.

511 After further investigation, we noticed that some reads classified as interrupted had a low ionic current at the end of the
512 read, while other reads had a similar ionic current level at the end of the read to fully-read reads (Fig. 6A). We divided
513 the interrupted reads into two classes based on their last 20 ionic current measurements: passing-like reads with high
514 end-mean, and blocking reads with low end-mean (< 50 pA). The passing-like reads exhibit a similar signal level at
515 the end of sequencing as fully-read reads, and the blocking reads exhibit a lower signal level. We analyzed the end
516 positions of these reads by extracting 25 bases upstream and downstream (50 in total) of the mapping end position in
517 the reference genome to study whether the reads typically ended at specific sequence motifs, which might be binding
518 sites of adducts. We also performed a motif discovery analysis on subsequences at the mapping end position and on
519 sequences that interact with the motor protein.

520 4.12 Motif affinity analysis

521 To estimate how well a k -mer matches a sequence motif, we applied AMA [46] to compute average motif affinity
522 scores. We did this for motifs detected with STREME on all possible k -mers, where $k = 9$, with the `--norc` option to
523 consider motif affinity only in the forward strand direction. In Figure 5A, a 9-mer was colored to match a motif if the
524 AMA score was more than 5. In the case where a 9-mer had high affinity to multiple motifs, we considered only the
525 motif with the highest affinity.

526 4.13 Analysis of DNA triplets at specific pore positions

527 We computed the effect of triplet-pore position pairs on quantile differences by iterating all seven triplet positions for a
528 9-mer by computing the mean quantile difference of all 9-mers that contain a triplet t at position i :

$$mean_qd_{t,i} = \frac{1}{4^6} \sum_{k \in K[i:i+3]=t} q_treated_k - q_control_k$$

529 where K is the set of all 9-mers.