# Chromatin is an ancient innovation conserved between Archaea and Eukarya

Ron Ammar[1], Dax Torti[1], Kyle Tsui[1,2], Marinella Gebbia[1,2], Tanja Durbic[1], Gary D. Bader[1,3], Guri Giaever[1,2] and Corey Nislow[1,4]

1 Department of Molecular Genetics, University of Toronto, Toronto, Ontario; Donnelly Centre for Cellular and Biomedical Research, University of Toronto, Toronto, ON
2 Department of Pharmaceutical Sciences, University of Toronto, Toronto, ON
3 Department of Computer Science, University of Toronto, Toronto, ON
4 Banting and Best Department of Medical Research, University of Toronto, Toronto, ON

## Abstract
The eukaryotic nucleosome is the fundamental unit of chromatin, comprising a protein octamer that wraps ~147bp of DNA and has essential roles in DNA compaction, replication and gene expression. Nucleosomes and chromatin have long been considered to be unique to eukaryotes, yet studies of select archaea have identified homologs of histone proteins that assemble into tetrameric nucleosomes. Here we report the first archaeal genome-wide nucleosome occupancy map, as observed in the halophile *Haloferax volcanii*. Nucleosome occupancy was compared with gene expression by compiling a comprehensive transcriptome of *Hfx. volcanii*. We found that archaeal transcripts possess hallmarks of eukaryotic chromatin structure: nucleosome-free regions at transcriptional start sites and conserved −1 and +1 promoter nucleosomes. Our observations demonstrate that histones and chromatin architecture evolved before the divergence of Archaea and Eukarya, suggesting that the fundamental role of chromatin in the regulation of gene expression is ancient.

ELIFE DIGEST

Single-celled microorganisms called archaea are one of the three domains of cellular life, along with bacteria and eukaryotes. Archaea are similar to bacteria in that they do not have nuclei, but genetically they have more in common with eukaryotes. Archaea are found in a wide range of habitats including the human colon, marshlands, the ocean and extreme environments such as hot springs and salt lakes.

It has been known since the 1990s that the DNA of archaea is wrapped around histones to form complexes that closely resemble the nucleosomes found in eukaryotes, albeit with four rather than eight histone subunits. Nucleosomes are the fundamental units of chromatin, the highly-ordered and compact structure that all the DNA in a cell is packed into. Now we know exactly how many nucleosomes are present in a given cell for some eukaryotes, notably yeast, and to a good approximation we know the position of each nucleosome during a variety of metabolic states and physiological conditions. We can also quantify the nucleosome occupancy, which is measure of the length of time that the nucleosomes spend in contact with the DNA: this is a critical piece of information because it determines the level of access that other proteins, including those that regulate gene expression, have to the DNA. These advances have been driven in large part by advances in technology, notably high-density microarrays for genome wide-studies of nucleosome occupancy, and massively parallel sequencing for direct nucleosome sequencing.

Ammar et al. have used these techniques to explore how the DNA of *Haloferax volcanii*, a species of archaea that thrives in the hyper-salty waters of the Dead Sea, is organized on a genome-wide basis. Despite some clear differences between the genomes of archaea and eukaryotes – for example, genomic DNA is typically circular in archaea and linear in eukaryotes – they found that the genome of *Hfx. volcanii* is organized into chromatin in a way that is remarkably similar to that seen in all eukaryotic genomes studied to date. This is surprising given that the chromatin in eukaryotes is confined to the nucleus, whereas there are no such constraints in archaea. In particular, Ammar et al. found that those regions of the DNA near the ends of

32 genes that mark where the transcription of the DNA into RNA should begin and end contain

33 have lower nucleosome occupancy than other regions. Moreover, the overall level of occupancy

34 in *Hfx. volcanii* was twice that of eukaryotes, which is what one would expect given that

35 nucleosomes in archaea contain half as many histone subunits as nucleosomes in eukaryotes.

36 Ammar et al. also confirmed that that the degree of nucleosome occupancy is correlated with

37 gene expression.

38

39 These two findings – the similarities between the chromatin in archaea and eukaryotes, and the

40 correlation between nucleosome occupancy and gene expression in archaea – raise an interesting

41 evolutionary possibility: the initial function of nucleosomes and chromatin formation might have

42 been for the regulation of gene expression rather than the packaging of DNA. This is consistent

43 with two decades of research that has shown that there is an extraordinarily complex relationship

44 between the structure of chromatin and the process of gene expression. It is possible, therefore,

45 that as the early eukaryotes evolved, nucleosomes and chromatin started to package DNA into

46 compact structures that, among other things, helped to prevent DNA damage, and that this

47 subsequently enabled the early eukaryotes to flourish.

48

## Introduction

Archaeal nucleosome core particles protect ~60 bp of DNA, approximately half that of eukaryotic nucleosomes, as demonstrated by the landmark work of Reeve and colleagues(Pereira et al., 1997). Comparing both eukaryotic and archaeal nucleosomes, the former is an octamer composed of heterodimers of histones H2A, H2B, H3 and H4 whereas the latter histones assemble from homologs of H3 and H4 proteins(Talbert and Henikoff, 2010, Pereira and Reeve, 1998). Archaeal histones can form both homodimers and heterodimers, as well as homotetramers, whereas eukaryotic histones contain hydrophobic dimerization surfaces that restrict assembly of the octamer from H2A-H2B and H3-H4 heterodimers(Sandman and Reeve, 2006, Talbert and Henikoff, 2010).

Using single-nucleotide resolution maps of archaeal nucleosome occupancy and gene expression, we demonstrate that the architecture of archaeal chromatin and the occupancy of its nucleosomes along transcription units are conserved. We constructed a nucleosome occupancy map of the halophilic archaeon *Haloferax volcanii*, a member of the phylum euryarchaeota, originally discovered in the highly saline sediment of the Dead Sea(Mullakhanbhai and Larsen, 1975). The genome of *Hfx. volcanii* has an average GC content of 65% and a total genome length of 4Mb(Hartman et al., 2010) composed of five circular genetic elements: a 2.8Mb main chromosome, three smaller chromosomes pHV1, pHV3 and pHV4 and the plasmid pHV2. It is highly polyploid with ~15 genome copies during exponential growth and ~10 during stationary phase(Breuert et al., 2006). The histone protein of *Hfx. volcanii*, hstA (HVO_0520), has a domain architecture containing two distinct histone fold domains in the same peptide that heterodimerize resembling that of the *Methanopyrus kandleri* histone (HMk)(Talbert and Henikoff, 2010, Marchler-Bauer et al., 2011, Geer et al., 2002).

72

**Results**

73

74       We cultured *Hfx. volcanii* in rich media containing 2M NaCl(Mullakhanbhai and Larsen,

75   1975). Genomic DNA was cross-linked and digested with micrococcal nuclease (MNase), with

76   cell disruption accomplished by bead-beating(Tsui et al., 2012). Nucleosome-bound cross-linked

77   genomic regions are protected from MNase digestion, in contrast to the linker DNA between

78   nucleosomes. Mononucleosome-sized (50-60bp) DNA fragments were gel purified and libraries

79   were sequenced on an Illumina HiSeq2000 (Fig. 1a). Sequence reads were aligned to the

80   published *Hfx. volcanii* DS2 genome(Hartman et al., 2010) to generate a genome-wide

81   nucleosome occupancy map. Controls included crosslinked DNA without MNase digestion as

82   well as MNase treated nucleosome-free genomic DNA.  The nucleosome occupancy data was

83   significantly different than the control MNase digest of deproteinized "naked" genomic DNA (r

84   = 0.071), indicating that the nucleosome map is unaffected by any potential MNase sequence

85   bias (Chung et al., 2010).

86       To determine nucleosome midpoints, we smoothed the occupancy data using a

87   symmetrical convolution sum with a Gaussian filter(Smith, 1997). Extrema were detected in the

88   smoothed signal, and maxima were defined as nucleosome midpoints. In the smoothed signal,

89   the mean peak-to-peak distance for the main chromosome was 68.5bp in genic regions and

90   76.1bp in non-genic regions. Genic regions were defined as the transcribed region plus 40bp (the

91   average promoter length based on Palmer and Daniels (1995)) upstream of the 5' end(Palmer and

92   Daniels, 1995). We observed a greater nucleosome density in *Hfx. volcanii* vs. all eukaryotes

93   likely due to the shorter length of DNA wrapped around the archaeal histone tetramer(Pereira et

94   al., 1997). Based on our data, the *Hfx. volcanii* genome has 14.2 nucleosomes/Kb compared to

5.2 nucleosomes/Kb in *Saccharomyces cerevisiae*. The resulting map reveals a periodic pattern similar to that seen in all eukaryotes examined to date; with protected regions appearing as peaks and linker regions as troughs. Sequence analysis of the entire nucleosome map showed that nucleosome midpoints were enriched with G/C nucleotides from 61.4% GC at the edge of the protected fragment to 74.6% GC at the midpoint (dyad). We found an increase of G/C nucleotides and a decrease in A/T nucleotides at the midpoint, as described recently for human cell lines (Fig. 1b,c)(Valouev et al., 2011). In contrast to previous studies in eukaryotes, we did not observe a periodicity in dinucleotide frequency relative to the nucleosome midpoint(Bailey et al., 2000, Satchwell et al., 1986, Albert et al., 2007).

We next investigated the relationship between nucleosome occupancy and gene expression. The existing genome annotation for *Haloferax* is derived almost exclusively from ORF predictions(Hartman et al., 2010). To augment these predictions, we used deep sequencing to create a high confidence transcriptome of the main chromosome of *Hfx. volcanii*. This map allowed us to define both 5'UTR lengths and transcriptional start sites (TSSs). Total RNA was extracted from *Hfx. volcanii* cells, repetitive RNA was partially depleted via duplex-specific nuclease (DSN) normalization followed by RNA-seq (see Methods)(Zhulidov et al., 2004). Transcript sequences were aligned, assembled and quantified using TopHat and the Genome Analysis Toolkit(Trapnell et al., 2009, McKenna et al., 2010) and transcript boundaries were further trimmed based on RNA-seq coverage information, as described previously(Wurtzel et al., 2010). The final set of transcripts were manually curated yielding 3059 transcriptional units in *Hfx. volcanii*, a number that is greater than observed previously in the comparable transcriptome of the sulfur-metabolizing archaeon *Sulfolobus solfataricus*(Wurtzel et al., 2010) but fewer than the 4073 predicted *Hfx. volcanii* genes. It is likely that under the rich media conditions used in

118     this study, not all genes are expressed. Specifically 75% of the predicted transcripts were

119     detectably expressed, and this fraction is consistent with observations obtained for yeast gene

120     expression in rich media(David et al., 2006). 32 novel transcripts (absent from the predicted

121     sequence annotation) were identified in the RNA-seq data. Most of these 32 transcripts lack

122     significant sequence homologs, and several were classified as transposases with paralogs in *Hfx.*

123     *volcanii* (Supplementary File 1). Notably, the gene that was most highly expressed in the

124     transcriptome (NTRANS_0004) was not previously annotated and contains a putative N-

125     Acyltransferase (NAT) superfamily domain. Homology searches revealed that this transcript

126     appears to be restricted to the genomes of other halophilic archaea (Altschul et al., 1990). The

127     architecture of this domain is homologous to chain A of the well-characterized histone

128     acetyltransferases Gcn5, Gna1, Hpa2 in *S. cerevisiae*, suggesting a possible role for this

129     transcript in regulating transcription via histone acetylation(Marchler-Bauer et al., 2011).

130     Additional acyltransferases with a similar architecture have been implicated in bacteriophage-

131     encoded DNA modifiers as well as cold and ethanol tolerance in yeast(Kaminska and Bujnicki,

132     2008, Du and Takagi, 2007). Thus, while post-translational modifications have not been

133     observed in archaeal histones (Forbes et al., 2004), our observation suggests that some

134     rudimentary control over chromatin accessibility may occur via the action of ancient NAT family

135     members. Furthermore acetyltransferase and deacetylase orthologs, which appear to have

136     enzymatic activity based on their sensitivity to the histone deacetylase (HDAC) inhibitor

137     trichostatin A have been identified in *Hfx. volcanii* (Altman-Price and Mevarech, 2009). In our

138     subsequent analysis, we focused on all genes we empirically determined to be expressed.

139        In eukaryotes, the TSS of the majority of expressed genes is characterized by a

140     nucleosome-depleted region (NDR)(Jiang and Pugh, 2009). This NDR is flanked by the well-

141    positioned −1 and the +1 nucleosomes. These regions direct RNA polymerase II to initiate

142    transcription and influence the binding of promoter regulatory elements(Jiang and Pugh, 2009).

143    This stereotypical pattern of nucleosome depletion at promoters and well-ordered nucleosomes

144    in gene bodies is found in all eukaryotes, including yeast, *Drosophila*, *A. thaliana* and humans.

145    Using the RNA-seq-derived transcripts for *Hfx. volcanii*, we computed the degree of aggregate

146    nucleosome occupancy for the 2343 transcripts on the main chromosome, and found that the

147    NDR and −1 and +1 nucleosomes are conserved in *Hfx. volcanii* (Fig. 2) suggesting that the

148    interplay between chromatin and transcription is conserved in archaeal promoters. We generated

149    nucleosome occupancy profiles for each transcript and clustered them hierarchically. Differential

150    nucleosome density was observed with profiles encompassing four to six nucleosomes in a

151    400bp DNA segment spanning 200bp on each side of the TSS (Fig. 2c). NDRs at transcription

152    termination sites (TTSs) are also observed, and similar to those found in eucaryotes (Lee et al.,

153    2007) they are less prominent than promoter NDRs in *Hfx. volcanii*.

154

155    **Discussion**

156         Our study establishes that nucleosome occupancy is conserved between archaea and

157    eukaryotes (Fig. 4). We further show that the nucleosomal protected fragments and NDRs are

158    shorter in archaea than in eukaryotes. Our findings are particularly noteworthy because *Hfx.*

159    *volcanii* likely resembles a deeply rooted ancestor that possessed eukaryotic genome architecture

160    hallmarks such as histones, as well as bacterial hallmarks such as the Shine-Dalgarno

161    sequence(Sartorius-Neef and Pfeifer, 2004). Archaeal histone tetramers likely resemble an

162    ancestral state of chromatin, as it has been observed that functional $(H3-H4)_2$ tetramers can be

163    formed *in vitro* from eukaryotic histones, and these tetramers are functional; they facilitate more

164    rapid transcription *in vitro* compared to native histone octamers(Puerta et al., 1993). The

165    observation that archaea contain (H3-H4)$_2$ tetramers is consistent with the proposal that

166    formation of the canonical eukaryotic nucleosome octamer begins with (H3-H4)$_2$ tetramer

167    assembly(Talbert and Henikoff, 2010).

168        Our study demonstrates that both histones and chromatin architecture evolved before the

169    divergence of Archaea and Eukarya, suggesting that the fundamental role of chromatin in the

170    regulation of gene expression is ancient. As well, owing to the small bacterial-sized archaeal

171    genome, we suggest that archaeal chromatin is not required for genome compaction. This leads

172    us to postulate that higher-order chromatin(Sajan and Hawkins, 2012) is a eukaryotic invention

173    and that archaeal chromatin is necessary but not sufficient for genome compaction. Furthermore

174    our observations provide a rich dataset that addresses the evolution of chromatin and its

175    fundamental role in the regulation of gene expression.

176

177    **Materials and Methods**

178    *Sample preparation. Haloferax volcanii* DS2 cells (obtained from the ATCC) were grown to

179    mid-log phase at 42°C in ATCC 974 Halobacterium medium supplemented with 2M NaCl. Cells

180    were fixed with 2% formaldehyde for 30 min then quenched with 125mM of glycine for 5 min.

181    An unfixed control sample was also prepared to serve as as a deproteinized, "naked" DNA

182    control, as described previously (Chung et al., 2010). Cells were pelleted and snap frozen prior to

183    MNase digestion and DNA extraction. Frozen cells were processed according to a modified

184    protocol from Rizzo *et al.*(Rizzo et al., 2011, Tsui et al., 2012). Samples were digested with

185    increasing concentrations of MNase and a no MNase control. After digestion, fragments 50-60bp

186    in length were size-selected using an Agilent Bioanalyzer High Sensitivity chip (Agilent, part#

187    5067-4626) and further processed for Illumina deep sequencing. This size-selection was critical,

188    as the formaldehyde crosslinking causes both histones as well as other DNA-binding proteins to

189    crosslink with bound DNA. Nucleosomal and genomic libraries were pooled equally according

190    to qPCR quantitation, and sequenced using v3 chemistry on one single-read HiSeq2000 lane

191    (50x8). Samples were demultiplexed using an 8bp index read at the end of read 1.

192

193    *Sequence read filtering and alignment.* Illumina sequencers require the ligation of an adapter

194    oligonucleotide to facilitate cluster formation on the flow cell. Because the library inserts were

195    short (~60bp), many sequence reads extended into the Illumina adapter sequences. The adapter

196    subsequences were computationally trimmed to ensure maximal read mapping. Then, using a

197    sequence quality cutoff of Phred20, reads were trimmed from both 5' and 3' ends to ensure

198    accurate mapping. These trimmed reads from control and MNase-treated genomic DNA were

199    aligned to the *Hfx. volcanii* DS2 genome using the Bowtie 2 gapped short read

200    aligner(Langmead and Salzberg, 2012). Sequence coverage was computed using the Genome

201    Analysis Toolkit (GATK) depth of coverage walker, which revealed the periodicity in the

202    occupancy data(DePristo et al., 2011).

203

204    *Nucleosome identification.* To detect nucleosome midpoint positions, sequence data were

205    Gaussian-smoothed as described previously by Shivaswamy *et al.* (2008) and Kaplan *et al.*

206    (2009)(Shivaswamy et al., 2008, Kaplan et al., 2009). This is appropriate because signals

207    generated by processes that are random, such as sequence coverage noise, usually have a

208    probability density function defined by a Gaussian distribution(Smith, 1997).

209    The Gaussian filter was defined as:

9

210  $$G(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{\left(\frac{-(x-\mu)^2}{2\sigma^2}\right)}$$

211  where $\mu$ is the mean of the distribution and $\sigma$ is the standard deviation.

212  A symmetrical convolution sum was applied with the following format:

213  $$y[i] = \sum_{j=-\frac{M}{2}}^{\frac{M}{2}} h[j] \cdot x[i-j]$$

214  where $M$ is an integer bandwidth, $y[j]$ is the output, $x[j]$ is the input and $h[j]$ is an $M$-point

215  function.

216  So, to smooth the coverage data, we applied the following convolution sum:

217  $$y[i] = \sum_{j=-\frac{M}{2}}^{\frac{M}{2}} G[j] \cdot x[i-j]$$

218  where $\sigma = \frac{M}{6}$. The interval length $M$ is constrained to $6\sigma$ because this encompasses 99.75% of the

219  Gaussian(Smith, 1997).

220  We also optimized nucleosome midpoint detection by convoluting a 2-pass simple moving

221  average (SMA) filter, but the Gaussian filter detected midpoints with greater resolution. Optimal

222  interval size for the Gaussian convolution sum, as determined by Pearson's correlation

223  coefficient with the raw data, was 27bp. For the 2-pass SMA it was 40bp for first-pass and 15bp

224  for second-pass.

225  Nucleosome occupancy was normalized genome-wide by transforming sequence coverage data

226  into binary-like data that existed in states of "occupied", "depleted" or transitioning between

227  those two states. This final occupancy map was used to define nucleosome positions.

228  Nucleosome occupancy profiles were clustered hierarchically by average linkage using Pearson's

229    correlation coefficient as the similarity metric in the Cluster 3.0 software package. Clusters were

230    visualized with Java Treeview (Fig. 2b,c).

231

232    *Transcript identification and genome annotation.* RNA was extracted with Trizol reagent

233    (Invitrogen, 15596-026), and DNase treated (Invitrogen, AM1907) according to manufacturer

234    specifications. A cDNA library was generated using 100ng of total RNA according to Illumina

235    TruSeq RNA Sample Prep protocol (Illumina, RS-122-2001) prior to duplex-specific nuclease

236    (DSN) treatment. 100ng of cDNA library was incubated in hybridization buffer (50mM HEPES,

237    500mM NaCl) for 2 minutes at 98°C, followed by 1 hour at 68°C. Ribosomal RNA (rRNA) was

238    not specifically depleted(He et al., 2010). Instead, we used duplex-specific nuclease (DSN)

239    normalization to remove recurrent RNA (rRNA, tRNA) from the total RNA sample, thereby

240    enriching mRNA(Zhulidov et al., 2004). Samples were immediately treated with 4 units of DSN

241    enzyme (Evrogen, EA001) in 1X DSN buffer and incubated for an additional 25 minutes at

242    68°C, prior to addition of stop solution, and purification with Ampure XP beads (Beckman

243    Coulter, A63881). RNA libraries were pooled equally according to qPCR quantitation, and

244    sequenced using v3 chemistry on a paired-end single HiSeq2000 lane (100x8x100). Samples

245    were demultiplexed using an 8bp index read at the end of read 1. Total RNA was sequenced at

246    extremely high coverage (2587× mean coverage) so that rRNA sequences (~77% of all sequence

247    reads) could be computationally excluded, as described by Wurtzel *et al.*(Wurtzel et al., 2010).

248    After quality score trimming (described earlier), sequence reads were aligned using

249    TopHat(Trapnell et al., 2009). The RNA-seq data displayed a great deal of overlap with the

250    predicted annotations(Hartman et al., 2010), with 92.1% of the existing annotations being

251    confirmed. Of the 4073 predicted annotations, 3751 were confirmed, and, of these, 744 were

252 merged with other transcripts to form longer transcripts. A heuristic approach was applied to

253 adjust the transcript 5' and 3' positions of the Hartman *et al.* predicted annotations based on the

254 boundaries of high RNA-seq coverage regions. This was vital as TSS accuracy is of great

255 importance for NDR identification (Fig. 5).

256 Because 85% of the *Haloferax* genome is predicted to be coding(Hartman et al., 2010), transcript

257 detection is complicated by transcript overlap. To overcome this, computationally identified

258 transcripts were manually curated yielding a total of 3059 expressed transcripts in *Hfx. volcanii*.

259 Of these, 32 transcripts are novel (Supplementary File 1). Of these transcripts, NTRANS_0004

260 was the most abundant transcript in the transcriptome, after the 6 rRNA genes. Homology data

261 was obtained using BLASTX with a BLOSUM45 matrix against the non-redundant protein

262 sequence database(Altschul et al., 1990). Conserved domains were identified using the

263 Conserved Domain Database(Marchler-Bauer et al., 2011). "Sequence data, nucleosome and

264 transcriptome maps and supplemental tables have been deposited to the Short Read Archive and

265 Dryad, as indicated in the datasets statement. Additionally this data is available at

266 http://chemogenomics.med.utoronto.ca/supplemental/chromatin/"

267

270
271 **References**

272 ALBERT, I., MAVRICH, T. N., TOMSHO, L. P., QI, J., ZANTON, S. J., SCHUSTER, S. C. & PUGH, B. F. 2007.
273       Translational and rotational settings of H2A.Z nucleosomes across the Saccharomyces cerevisiae genome.
274       *Nature,* 446**,** 572-6.

275 ALTMAN-PRICE, N. & MEVARECH, M. 2009. Genetic evidence for the importance of protein acetylation and
276       protein deacetylation in the halophilic archaeon Haloferax volcanii. *Journal of bacteriology,* 191**,** 1610-7.

277 ALTSCHUL, S. F., GISH, W., MILLER, W., MYERS, E. W. & LIPMAN, D. J. 1990. Basic local alignment search
278       tool. *Journal of molecular biology,* 215**,** 403-10.

279 BAILEY, K. A., PEREIRA, S. L., WIDOM, J. & REEVE, J. N. 2000. Archaeal histone selection of nucleosome
280         positioning sequences and the procaryotic origin of histone-dependent genome evolution. *Journal of*
281         *molecular biology,* 303**,** 25-34.

282 BREUERT, S., ALLERS, T., SPOHN, G. & SOPPA, J. 2006. Regulated polyploidy in halophilic archaea. *PloS one,*
283         1**,** e92.

284 CHANG, G. S., NOEGEL, A. A., MAVRICH, T. N., MULLER, R., TOMSHO, L. P., WARD, E., FELDER, M.,
285         JIANG, C., EICHINGER, L., GLOCKNER, G., SCHUSTER, S. C. & PUGH, B. F. 2012. Unusual
286         combinatorial involvement of poly-A/T tracts in organizing genes and chromatin in Dictyostelium. *Genome*
287         *research*.

288 CHUNG, H. R., DUNKEL, I., HEISE, F., LINKE, C., KROBITSCH, S., EHRENHOFER-MURRAY, A. E.,
289         SPERLING, S. R. & VINGRON, M. 2010. The effect of micrococcal nuclease digestion on nucleosome
290         positioning data. *PloS one,* 5**,** e15754.

291 DAVID, L., HUBER, W., GRANOVSKAIA, M., TOEDLING, J., PALM, C. J., BOFKIN, L., JONES, T., DAVIS,
292         R. W. & STEINMETZ, L. M. 2006. A high-resolution map of transcription in the yeast genome. *Proc Natl*
293         *Acad Sci U S A,* 103**,** 5320-5.

294 DEPRISTO, M. A., BANKS, E., POPLIN, R., GARIMELLA, K. V., MAGUIRE, J. R., HARTL, C.,
295         PHILIPPAKIS, A. A., DEL ANGEL, G., RIVAS, M. A., HANNA, M., MCKENNA, A., FENNELL, T. J.,
296         KERNYTSKY, A. M., SIVACHENKO, A. Y., CIBULSKIS, K., GABRIEL, S. B., ALTSHULER, D. &
297         DALY, M. J. 2011. A framework for variation discovery and genotyping using next-generation DNA
298         sequencing data. *Nature genetics,* 43**,** 491-8.

299 DU, X. & TAKAGI, H. 2007. N-Acetyltransferase Mpr1 confers ethanol tolerance on Saccharomyces cerevisiae by
300         reducing reactive oxygen species. *Applied microbiology and biotechnology,* 75**,** 1343-51.

301 FIUME, M., WILLIAMS, V., BROOK, A. & BRUDNO, M. 2010. Savant: genome browser for high-throughput
302         sequencing data. *Bioinformatics,* 26**,** 1938-44.

303 FORBES, A. J., PATRIE, S. M., TAYLOR, G. K., KIM, Y. B., JIANG, L. & KELLEHER, N. L. 2004. Targeted
304         analysis and discovery of posttranslational modifications in proteins from methanogenic archaea by top-
305         down MS. *Proceedings of the National Academy of Sciences of the United States of America,* 101**,** 2678-83.

306 GEER, L. Y., DOMRACHEV, M., LIPMAN, D. J. & BRYANT, S. H. 2002. CDART: protein homology by domain
307         architecture. *Genome research,* 12**,** 1619-23.

308 HARTMAN, A. L., NORAIS, C., BADGER, J. H., DELMAS, S., HALDENBY, S., MADUPU, R., ROBINSON, J.,
309         KHOURI, H., REN, Q., LOWE, T. M., MAUPIN-FURLOW, J., POHLSCHRODER, M., DANIELS, C.,
310         PFEIFFER, F., ALLERS, T. & EISEN, J. A. 2010. The complete genome sequence of Haloferax volcanii
311         DS2, a model archaeon. *PloS one,* 5**,** e9605.

312 HE, S., WURTZEL, O., SINGH, K., FROULA, J. L., YILMAZ, S., TRINGE, S. G., WANG, Z., CHEN, F.,
313         LINDQUIST, E. A., SOREK, R. & HUGENHOLTZ, P. 2010. Validation of two ribosomal RNA removal
314         methods for microbial metatranscriptomics. *Nature methods,* 7**,** 807-12.

315 JIANG, C. & PUGH, B. F. 2009. Nucleosome positioning and gene regulation: advances through genomics. *Nature*
316         *reviews. Genetics,* 10**,** 161-72.

317 KAMINSKA, K. H. & BUJNICKI, J. M. 2008. Bacteriophage Mu Mom protein responsible for DNA modification
318         is a new member of the acyltransferase superfamily. *Cell cycle,* 7**,** 120-1.

319 KAPLAN, N., MOORE, I. K., FONDUFE-MITTENDORF, Y., GOSSETT, A. J., TILLO, D., FIELD, Y.,
320         LEPROUST, E. M., HUGHES, T. R., LIEB, J. D., WIDOM, J. & SEGAL, E. 2009. The DNA-encoded
321         nucleosome organization of a eukaryotic genome. *Nature,* 458**,** 362-6.

322 LANGMEAD, B. & SALZBERG, S. L. 2012. Fast gapped-read alignment with Bowtie 2. *Nature methods,* 9**,** 357-9.

323 MARCHLER-BAUER, A., LU, S., ANDERSON, J. B., CHITSAZ, F., DERBYSHIRE, M. K., DEWEESE-SCOTT,
324         C., FONG, J. H., GEER, L. Y., GEER, R. C., GONZALES, N. R., GWADZ, M., HURWITZ, D. I.,
325         JACKSON, J. D., KE, Z., LANCZYCKI, C. J., LU, F., MARCHLER, G. H., MULLOKANDOV, M.,
326         OMELCHENKO, M. V., ROBERTSON, C. L., SONG, J. S., THANKI, N., YAMASHITA, R. A.,

327  ZHANG, D., ZHANG, N., ZHENG, C. & BRYANT, S. H. 2011. CDD: a Conserved Domain Database for
328      the functional annotation of proteins. *Nucleic acids research,* 39, D225-9.

329  MCKENNA, A., HANNA, M., BANKS, E., SIVACHENKO, A., CIBULSKIS, K., KERNYTSKY, A.,
330      GARIMELLA, K., ALTSHULER, D., GABRIEL, S., DALY, M. & DEPRISTO, M. A. 2010. The Genome
331      Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome*
332      *research,* 20**,** 1297-303.

333  MULLAKHANBHAI, M. F. & LARSEN, H. 1975. Halobacterium volcanii spec. nov., a Dead Sea halobacterium
334      with a moderate salt requirement. *Archives of microbiology,* 104**,** 207-14.

335  PALMER, J. R. & DANIELS, C. J. 1995. In vivo definition of an archaeal promoter. *Journal of bacteriology,* 177**,**
336      1844-9.

337  PEREIRA, S. L., GRAYLING, R. A., LURZ, R. & REEVE, J. N. 1997. Archaeal nucleosomes. *Proceedings of the*
338      *National Academy of Sciences of the United States of America,* 94**,** 12633-7.

339  PEREIRA, S. L. & REEVE, J. N. 1998. Histones and nucleosomes in Archaea and Eukarya: a comparative analysis.
340      *Extremophiles : life under extreme conditions,* 2**,** 141-8.

341  PUERTA, C., HERNANDEZ, F., GUTIERREZ, C., PINEIRO, M., LOPEZ-ALARCON, L. & PALACIAN, E.
342      1993. Efficient transcription of a DNA template associated with histone (H3.H4)2 tetramers. *The Journal of*
343      *biological chemistry,* 268**,** 26663-7.

344  RIZZO, J. M., MIECZKOWSKI, P. A. & BUCK, M. J. 2011. Tup1 stabilizes promoter nucleosome positioning and
345      occupancy at transcriptionally plastic genes. *Nucleic acids research,* 39**,** 8803-19.

346  SAJAN, S. A. & HAWKINS, R. D. 2012. Methods for Identifying Higher-Order Chromatin Structure. *Annual*
347      *review of genomics and human genetics*.

348  SANDMAN, K. & REEVE, J. N. 2006. Archaeal histones and the origin of the histone fold. *Current opinion in*
349      *microbiology,* 9**,** 520-5.

350  SARTORIUS-NEEF, S. & PFEIFER, F. 2004. In vivo studies on putative Shine-Dalgarno sequences of the
351      halophilic archaeon Halobacterium salinarum. *Molecular microbiology,* 51**,** 579-88.

352  SATCHWELL, S. C., DREW, H. R. & TRAVERS, A. A. 1986. Sequence periodicities in chicken nucleosome core
353      DNA. *Journal of molecular biology,* 191**,** 659-75.

354  SHIVASWAMY, S., BHINGE, A., ZHAO, Y., JONES, S., HIRST, M. & IYER, V. R. 2008. Dynamic remodeling
355      of individual nucleosomes across a eukaryotic genome in response to transcriptional perturbation. *PLoS*
356      *biology,* 6**,** e65.

357  SMITH, S. W. 1997. *The scientist and engineer's guide to digital signal processing,* San Diego, Calif., California
358      Technical Pub.

359  TALBERT, P. B. & HENIKOFF, S. 2010. Histone variants--ancient wrap artists of the epigenome. *Nature reviews.*
360      *Molecular cell biology,* 11**,** 264-75.

361  TRAPNELL, C., PACHTER, L. & SALZBERG, S. L. 2009. TopHat: discovering splice junctions with RNA-Seq.
362      *Bioinformatics,* 25**,** 1105-11.

363  TSUI, K., DURBIC, T., GEBBIA, M. & NISLOW, C. 2012. Genomic approaches for determining nucleosome
364      occupancy in yeast. *Methods in molecular biology,* 833**,** 389-411.

365  VALOUEV, A., JOHNSON, S. M., BOYD, S. D., SMITH, C. L., FIRE, A. Z. & SIDOW, A. 2011. Determinants of
366      nucleosome organization in primary human cells. *Nature,* 474**,** 516-20.

367  WURTZEL, O., SAPRA, R., CHEN, F., ZHU, Y., SIMMONS, B. A. & SOREK, R. 2010. A single-base resolution
368      map of an archaeal transcriptome. *Genome research,* 20**,** 133-41.

369  ZHULIDOV, P. A., BOGDANOVA, E. A., SHCHEGLOV, A. S., VAGNER, L. L., KHASPEKOV, G. L.,
370      KOZHEMYAKO, V. B., MATZ, M. V., MELESHKEVITCH, E., MOROZ, L. L., LUKYANOV, S. A. &
371      SHAGIN, D. A. 2004. Simple cDNA normalization using kamchatka crab duplex-specific nuclease.
372      *Nucleic acids research,* 32**,** e37.

373

**Figure Titles and Legends**

**Fig. 1.**

**Micrococcal nuclease digestion produces nucleosomal fragments from crosslinked *Hfx.***

***volcanii* chromatin. (A)** Formaldehyde cross-linked chromatin was subjected to MNase

digestion with increasing amounts on microccocal nuclease (from 1 unit to 5 units). De-

crosslinked DNAs were separated on a 3% agarose gel and **~60bp and ~120bp mono- and di-**

**nucleosomes were observed**. Markers (M) indicate * 50bp and ** 150bp. (**B**) The counts of AA,

AT, TA, TT or CC, CG, GC, GG dinucleotides are reported at each position showing an

enrichment of G/C nucleotides and a depletion of A/T nucleotides at the dyad relative to the end

points of the protected fragment. This differs from the observation of Bailey *et al.* (2000), where

GC, AA and TA dinucleotides were repeated at ~10bp intervals in recombinant archaeal histone

B from *Methanothermus fervidus* (rHMfB)(Bailey et al., 2000). (**C**) The sequence logo of a

nucleosome-binding site in *Hfx. volcanii* centered at the nucleosome midpoint. There is a

significant GC enrichment towards the nucleosome midpoint. This is exhibited using both bit

score and probability measures.

390

**Fig. 2**. **Nucleosome occupancy in *Haloferax volcanii*. (A)** Degree of normalized nucleosome

occupancy in aggregate for the main chromosome. As observed in eukaryotes, there is a

prominent nucleosome-depleted region (NDR) at the transcriptional start site (TSS) preceded by

a −1 nucleosome and followed by a +1 nucleosome, demonstrating that promoter genome

architecture is conserved between archaea and eukaryotes. (**B**) Hierarchical clustergram for the

2343 expressed transcripts on the main *Haloferax* chromosome. Green represents nucleosome-

depleted regions and red represents occupied regions. (**C**) The clustered heatmap was subdivided

398   into the largest 6 subclades, and differential density of nucleosomes can be observed with

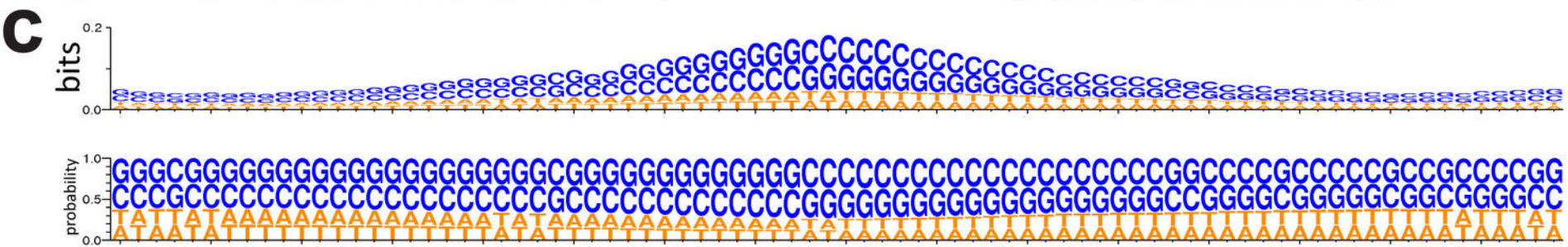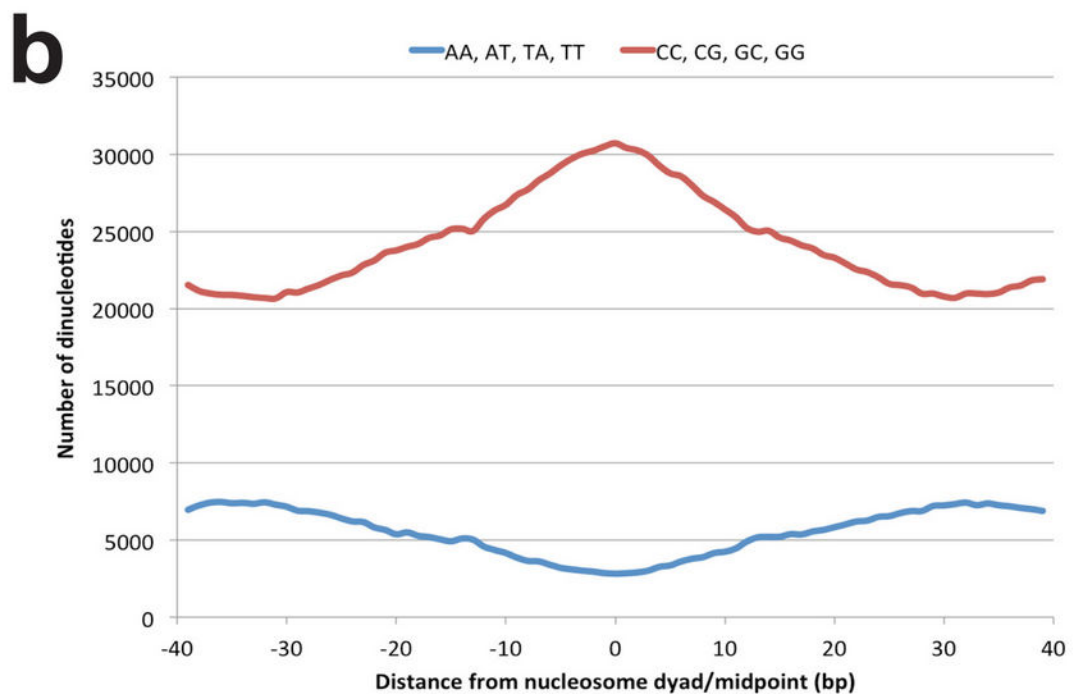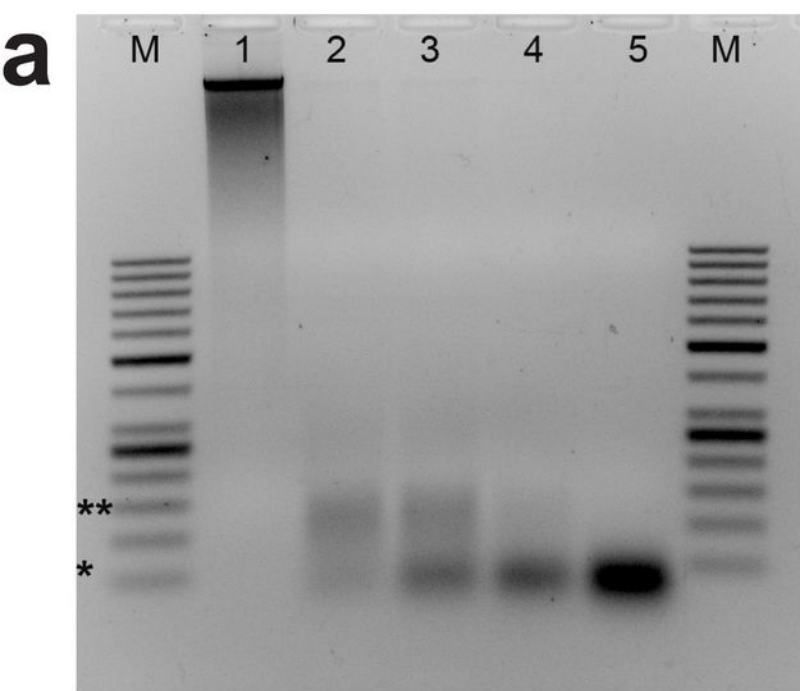399   occupancy profile clusters containing between 4 to 6 nucleosomes.

400

401   **Fig. 3. Nucleosome-depleted regions at the 3' end of transcripts.** As observed in eukaryotes,

402   NDRs are also found at the transcriptional termination sites in *Hfx. volcanii*. Both 5' and 3' end

403   profiles are overlaid in this figure for comparison. The 5' NDR is, on average, more depleted and

404   longer.

405

406   **Fig. 4. Chromatin architecture is conserved at the 5' end of transcripts across eukaryotes**

407   **and archaea.** Due to the smaller size of archaeal nucleosome DNA, the occupancy has a shorter

408   periodicity. Figure adapted with permission from Chang *et al*.(Chang et al., 2012).
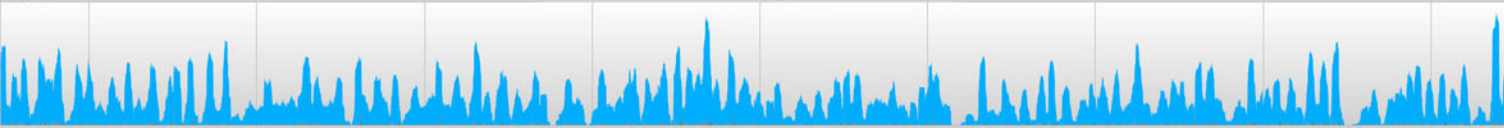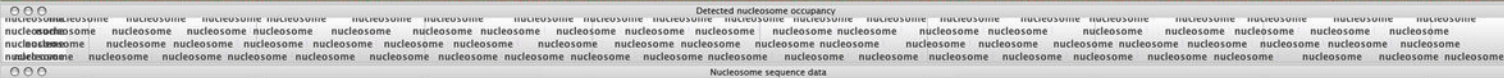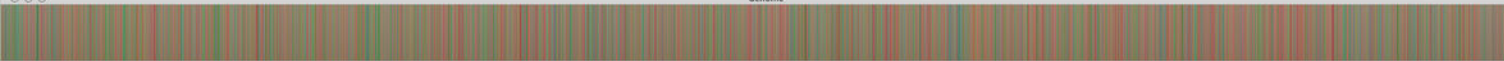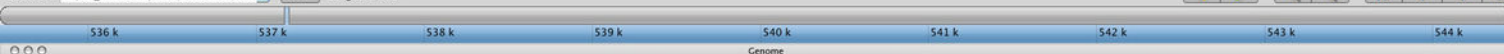
409

410   **Fig. 5. Sample screenshot of all data tracks loaded into the Savant genome browser (Fiume**

411   **et al., 2010).** The nucleosome sequence data is displayed, and the periodicity reflects protected

412   and unprotected fragments after MNase digestion (magnitude of peak is not considered). Peaks

413   represent nucleosome midpoints, which were detected and marked. Below are the corresponding

414   RNA-seq and curated gene tracks. In this screenshot, one can observe seven entire ORFs in line

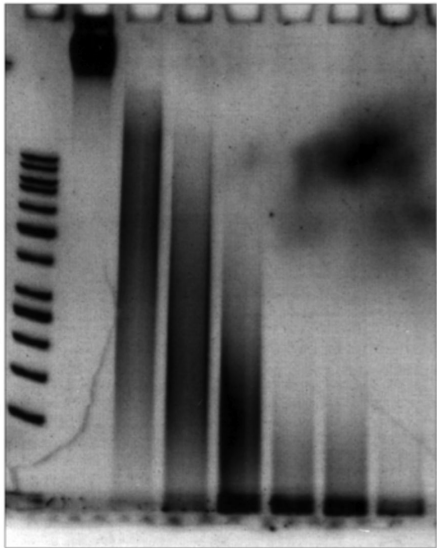415   with their NDRs and −1 and +1 nucleosomes.

*S. pombe*

TSS

+1
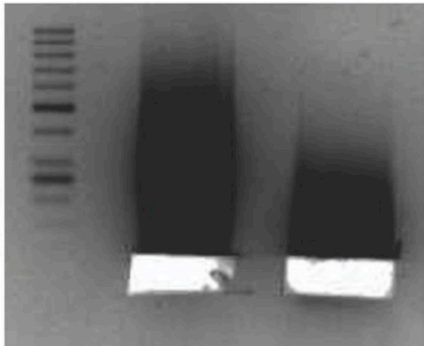
*S. cerevisiae*

+1

*C. elegans*

+1

*D. melanogaster*

+1

*Human*

+1

*D. discoideum*

+1

*A. thaliana*

+1

*Hfx. volcanii*

+1

-400    0    400

Distance from TSS

Location: ref|NC_013967.1|: 535,473 - 544,473    Go    Length: 9,001

536 k    537 k    538 k    539 k    540 k    541 k    542 k    543 k    544 k

Genome

Detected nucleosome occupancy

nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome nucleosome

Nucleosome sequence data

RNA-seq data

Curated genes

CDS_HVO_0605    CDS_HVO_0607    CDS_HVO_0609    CDS_HVO_0613
CDS_HVO_0604    CDS_HVO_0606    CDS_HVO_0608    CDS_HVO_0610+    CDS_HVO_0612

sample 2