











# The 4D Nucleome Data Portal as a resource for searching and visualizing curated nucleomics data

Sarah B. Reiff <sup>1</sup>, Andrew J. Schroeder <sup>1</sup>, Koray Kırılı<sup>1</sup>, Andrea Cosolo <sup>1</sup>, Clara Bakker <sup>1</sup>, Luisa Mercado<sup>1</sup>, Soohyun Lee<sup>1</sup>, Alexander D. Veit<sup>1</sup>, Alexander K. Balashov<sup>1</sup>, Carl Vitzthum<sup>1</sup>, William Ronchetti<sup>1</sup>, Kent M. Pitman <sup>1</sup>, Jeremy Johnson<sup>1</sup>, Shannon R. Ehmsen<sup>1</sup>, Peter Kerpedjiev<sup>1</sup>, Nezar Abdennur<sup>2</sup>, Maxim Imakaev <sup>2</sup>, Serkan Utku Öztürk<sup>3</sup>, Uğur Çamoğlu<sup>3</sup>, Leonid A. Mirny <sup>2,4</sup>, Nils Gehlenborg <sup>1</sup>, Burak H. Alver <sup>1</sup> & Peter J. Park <sup>1</sup>✉

The 4D Nucleome (4DN) Network aims to elucidate the complex structure and organization of chromosomes in the nucleus and the impact of their disruption in disease biology. We present the 4DN Data Portal (<https://data.4dnucleome.org/>), a repository for datasets generated in the 4DN network and relevant external datasets. Datasets were generated with a wide range of experiments, including chromosome conformation capture assays such as Hi-C and other innovative sequencing and microscopy-based assays probing chromosome architecture. All together, the 4DN data portal hosts more than 1800 experiment sets and 36000 files. Results of sequencing-based assays from different laboratories are uniformly processed and quality-controlled. The portal interface allows easy browsing, filtering, and bulk downloads, and the integrated HiGlass genome browser allows interactive visualization and comparison of multiple datasets. The 4DN data portal represents a primary resource for chromosome contact and other nuclear architecture data for the scientific community.

<sup>1</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02115, USA. <sup>2</sup>Institute for Medical Engineering and Science, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. <sup>3</sup>Karya SMD Software Solutions, İzmir 35040, Turkey. <sup>4</sup>Department of Physics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. ✉email: [peter\\_park@hms.harvard.edu](mailto:peter_park@hms.harvard.edu)

The 4D Nucleome (4DN) Network<sup>1</sup> is an NIH Common Fund project that started in 2015 with the overarching goal of elucidating the three-dimensional organization of chromosomes in the nucleus across different cell types and cell cycle stages, and to understand how perturbation of this structure can impact human health in cases of cancer and other diseases. The first phase (2015–2020) of the consortium included five projects on nucleomics technology development, nine projects on the development of imaging tools, six projects studying nuclear bodies and compartments, six centers studying nuclear organization, and 10 additional collaborating projects, involving ~100 laboratories in total. The second phase of the consortium (2020–2025) includes eight projects studying real-time chromatin dynamics, 16 projects examining the role of nuclear biology in human health and disease, and four centers developing methods for data integration, modeling, and visualization. The new projects bring the total of 4DN network laboratories to over 150 and promise the generation of even more nucleomics data, particularly with relevance to human health, as well as more advanced modes of data integration across disparate assay types.

The 4DN Data Coordination and Integration Center (DCIC) has created a web portal that serves as a repository for data generated by the 4DN Network members. It has also imported external datasets that are widely used by the community. Implemented fully on the Amazon Web Services (AWS) platform with the latest technologies, the portal has been engineered to provide utility to the broader community of nuclear biology researchers. It enables easy searching and browsing of the data and, importantly, the associated metadata, thus allowing for increased reproducibility both at the analytical and experimental levels.

The portal was designed to accommodate data generated from both genomics and microscopy experiments. A large portion of the genomics data comes from chromosome conformation capture (3C) assays such as Hi-C<sup>2</sup>. Additional data derived from several other genomic assay types that probe chromosome conformation or other aspects of nuclear structure and function, such as replication timing, chromatin accessibility, or spatial proximity of DNA to other cellular components. These assays include Repli-seq<sup>3</sup>, DamID<sup>4</sup>, CUT&RUN<sup>5</sup>, and SPRITE<sup>6</sup>. We have incorporated elements specific to imaging modalities into our data model as well, and provide visualization capabilities for some imaging data types. Many current datasets use imaging as a tool for measuring distances between genomic locations or between a genomic region and a sub-nuclear structure, and imaging is often used as a cross-validation technique for the results of genomic assays. Current imaging data types available include standard fluorescent in situ hybridization (FISH) that targets a DNA or RNA sequence, multiplexed FISH that can target many loci in a single fixed sample<sup>7</sup>, and dynamic single-particle tracking.

In addition to making data easily accessible to users, the 4DN DCIC has aimed to develop and run standardized bioinformatics pipelines on submitted raw data to generate consistent and comparable results. These can then be explored in an integrated HiGlass browser that supports visualization of 2D contact maps and linear 1D tracks such as gene annotations or CUT&RUN peaks<sup>8</sup>. The portal also supports the download of raw and processed files for users that want to use the files locally.

Here we describe the 4DN Data Portal, a repository for genomic and microscopy nuclear architecture datasets. We discuss data that can be found in the portal, interactive visualization capabilities for these datasets in the portal, and utility to users outside the 4DN network.

## Results

**4DN Data at the data portal.** The majority of datasets in the portal are submitted by the 4DN researchers who performed the studies. A subset of datasets on the 4DN data portal was generated outside the consortium but was curated by 4DN curators because they were widely used in the community; these data were uploaded from other public repositories such as NCBI's Gene Expression Omnibus (GEO) and Sequence Read Archive (SRA). External users interested in submitting their data or those wishing to suggest a dataset for inclusion in the portal may contact our support desk ([support@4dnucleome.org](mailto:support@4dnucleome.org)).

Within the data model of the 4DN data portal, the top-level item is called an Experiment Set, and represents a set of replicate experiments, all performed under the same experimental conditions; as such, quantities of data will be discussed below in terms of experiment sets. 4DN data generators are strongly encouraged to include multiple biological replicates for every experiment, and currently, over 800 experiment sets in the portal have met this standard.

**Genomics datasets.** For genomics assays, Hi-C and its variants make up the largest proportion of experiments in the 4DN data portal. The original Hi-C protocol was reported in 2009<sup>2</sup>, as a chromosome conformation capture (3C) method that assays pairwise interactions in DNA across the whole genome by making use of high-throughput sequencing. Since then, many variants of Hi-C have been described. The first Hi-C experiments used 1% formaldehyde cross-linking, a six-cutter restriction enzyme, and digestion and ligation steps were performed in lysed nuclei in a dilute solution. Later, an in situ version involved performing the digestion and ligation steps in intact nuclei. Other types of Hi-C assays involve using different cross-linking chemistries and chromatin fragmentation methods (Micro-C, DNase Hi-C, Hi-C 3.0)<sup>9–11</sup>, and single-cell and single nucleus versions of Hi-C have also been implemented<sup>12–14</sup>. In total, the 4DN data portal has 610 public Hi-C experiment sets, spanning nine subtypes of Hi-C assays (Table 1). Some of these datasets are particularly useful because of their high sequencing depth. Over 40 Hi-C experiment sets on the portal have over 1 billion read pairs after filtering, and a table of these is available on the data portal resource pages (<https://data.4dnucleome.org/resources/data-collections/high-resolution-hic-datasets>).

The same year that Hi-C was reported, the first ChIA-PET (Chromatin Interaction Analysis by Paired-End Tag sequencing) experiments were published. Similar to Hi-C in design, ChIA-PET uses chromatin immunoprecipitation to enrich DNA-protein complexes and then employs proximity ligation and sequencing, enabling the discovery of chromatin interactions mediated by a target protein<sup>15,16</sup>. Together with ChIA-Drop, a droplet-based assay similar to ChIA-PET<sup>17</sup>, and PLAC-seq, another similar assay in which proximity ligation is performed prior to immunoprecipitation<sup>18</sup>, IP-based 3C assays comprise 34 public experiment sets on the 4DN data portal (Table 1).

Recently, a great deal of technology has been developed by members of the 4DN community, resulting in new innovative assays for probing aspects of nuclear DNA structure. To give a few examples, SPRITE is a method for identifying multi-way interactions between distal genomic regions<sup>6</sup>; MARGI is a method for mapping RNA-chromatin interactions<sup>19–21</sup>; CUT&RUN is a DNA-binding assay generating data similar to ChIP-seq but which produces much lower background<sup>5</sup>; and TSA-seq is being used to analyze the proximity of nuclear structures to genomic regions on a genome-wide scale<sup>22</sup>. All the genomic assays currently on the 4DN data portal and available to the public are shown in (Table 1).

**Table 1 Genomic assay types in the 4D Nucleome Data Portal.**

| Experiment type                  | No. of public experiment sets |
|----------------------------------|-------------------------------|
| <b>Hi-C</b>                      | <b>610</b>                    |
| in situ Hi-C                     | 335                           |
| Dilution Hi-C                    | 118                           |
| DNase Hi-C                       | 21                            |
| Micro-C                          | 26                            |
| Single-cell Hi-C                 | 11                            |
| Single nucleus Hi-C              | 17                            |
| sci-Hi-C <sup>13</sup>           | 28                            |
| Capture Hi-C <sup>53</sup>       | 40                            |
| TCC                              | 14                            |
| <b>IP-based 3C assays</b>        | <b>34</b>                     |
| ChIA-PET                         | 4                             |
| in situ ChIA-PET                 | 10                            |
| ChIA-Drop <sup>17</sup>          | 2                             |
| PLAC-seq <sup>18</sup>           | 18                            |
| <b>DNA-binding assays</b>        | <b>202</b>                    |
| ChIP-seq                         | 141                           |
| CUT&RUN                          | 61                            |
| CUT&Tag <sup>54</sup>            | 2                             |
| <b>Other sequencing assays</b>   | <b>421</b>                    |
| Repli-seq <sup>55</sup>          | 138                           |
| SPRITE <sup>6</sup>              | 3                             |
| DamID                            | 66                            |
| ATAC-seq                         | 21                            |
| RNA-seq                          | 90                            |
| TRIP <sup>56</sup>               | 7                             |
| NAD-seq <sup>57</sup>            | 8                             |
| TSA-seq <sup>22</sup>            | 67                            |
| MARGI <sup>19</sup>              | 6                             |
| GAM <sup>58</sup>                | 6                             |
| RE-seq(DpnII-seq <sup>59</sup> ) | 11                            |
| Bru-seq <sup>60</sup>            | 1                             |
| MC-3C <sup>61</sup>              | 1                             |
| <b>Total</b>                     | <b>1273</b>                   |

Bold rows indicate categories of genomic assays and their subtotal counts.

**Genomics pipelines and QC.** For genomics datasets, data are submitted as raw fastq files. All data are aligned to the genome references hg38 and mm10 for analyses of human and mouse samples, respectively. For Hi-C experiments, the general analysis pipeline involves (1) aligning reads to a reference genome; (2) filtering aligned reads; (3) combining replicates for a single list of contact pairs; and (4) matrix aggregation and normalization. At each of these steps, different tools and parameters are employed by independent laboratories. Although optimal tools/parameters may depend on the type of downstream analysis performed, a uniform processing pipeline was necessary to ensure consistency and compatibility. After extensive discussions and testing by the members of the 4DN Analysis Working Group, the final version of the Hi-C processing included the following steps: (1) alignment by BWA MEM<sup>23</sup> with the -SP5M option to ensure that paired reads are aligned independently but the results are formatted properly as paired-end data and that the 5' portion of a chimeric alignment is reported as a primary soft-clipped alignment, (2) sorting and filtering the reads using pairtools<sup>24</sup> and (3) aggregating filtered reads into a contact matrix and normalizing it. The analysis also includes quality control steps using Fastqc<sup>25</sup> and Pairsqc<sup>26</sup>. This pipeline outputs a multi-resolution contact matrix at the 4DN standard resolutions of 1 kb, 2 kb, 5 kb, 10 kb, 25 kb, 50 kb, 100 kb, 250 kb, 500 kb, 1 Mb, 2.5 Mb, 5 Mb and 10 Mb. This matrix is generated using the cooler software and the .mcool file format<sup>27</sup>, which is compatible with the HiGlass interactive 2D

genome browser<sup>8</sup> and can be visualized on the portal. An additional contact matrix is also generated in .hic format, which is compatible with the Juicebox 2D interactive genome browser<sup>28</sup>. More details of the pipeline can be found in the portal's resource pages at <https://data.4dnucleome.org/resources/data-analysis/hic-processing-pipeline>.

Genomic contact matrices generally show evidence of genomic compartments<sup>2</sup> as well as local regions of enriched intra-compartment contacts, known as topologically-associated domains (TADs)<sup>29,30</sup>. There is a great deal of interest in exploring the nature and dynamics of TADs and sub-TADs, domains nested within others. Identification of such domains is difficult due to the lack of validated sets, and detection algorithms are continually being refined<sup>31-33</sup>. At this time, the 4DN Data Portal runs two domain identification workflows on contact matrices to report compartments and TAD boundaries, using the cooltools software<sup>34</sup>. The first workflow uses an eigenvector decomposition of the matrix to call active (A) and inactive (B) compartments. For cis contacts, eigenvalues are calculated by the cooltools command `cooler_cis_eig` with default values of `n_eigs=3`, `phasing_track_col='GC'`, `ignore_diags=None`, `clip_percentile=99.9`, and `sort_metric='var_explained'`. For trans contacts, the `cooler_trans_eig` command is used with the default parameters `n_eigs=3`, `partition=None`, `phasing_track_col='GC'`, and `sort_metric='var_explained'`. The second workflow computes insulation scores along the diagonal of the matrix, based on average interaction frequencies crossing over each genomic bin, and prominent dips in this score indicate boundaries between domains<sup>35</sup>. The insulation table is calculated with the cooltools `calculate_insulation_score` command with default parameter `window=100000`, followed by the `find_boundaries` command with default parameters `pixels_frac=0.66` and `cutoff=2`. As the nature of TADs and domains is an area of active investigation, the 4DN DCIC chooses to report boundaries based on insulation scores to provide results that might be further utilized in developing domain-detection algorithms. Results of both workflows are accessible as bed files as well as bigwig files which can be visualized in HiGlass on the data portal. More details on the domain-detection pipelines can be found at [https://data.4dnucleome.org/resources/data-analysis/insulation\\_compartment\\_scores](https://data.4dnucleome.org/resources/data-analysis/insulation_compartment_scores).

Processing pipelines are run on other types of high-throughput sequencing assays at the portal as well, including Repli-seq, CUT&RUN, and MARGI (Table 2). The pipeline for Repli-seq involves (1) trimming reads with cutadapt and aligning to a reference genome with bwa-mem<sup>23</sup>; (2) filtering valid alignments with samtools<sup>36</sup>; and (3) binning and aggregating for 5kb windows with bedtools<sup>37</sup>. The final output is provided in gzipped bedgraph and bigWig formats.

The CUT&RUN pipeline on the data portal processes paired-end reads in three main steps: (1) reads are processed with Trimmomatic for quality filtering and adapter trimming, and aligned to a reference genome with bowtie2<sup>38</sup>; (2) duplicate reads are filtered out using Picard<sup>39</sup> and samtools<sup>36</sup>, and converted into .bed format with bedtools<sup>37</sup>; and (3) peaks are called with SEACR<sup>40</sup>. The final outputs of the pipeline include a peaks file and a bigWig track. The final bigWig files of the CUT&RUN and Repli-seq pipelines can both be visualized on the portal as 1D tracks in the HiGlass browser. The MARGI pipeline is similar to the Hi-C pipeline, and is adapted from the original pipeline written by the iMARGI creators<sup>21</sup>.

For ATAC-seq, ChIP-seq, and RNA-seq data, pipelines from the ENCODE Data Coordination Center<sup>41</sup> have been adapted for the 4DN platform. For ATAC-seq and ChIP-seq, the final output of these pipelines includes a bigwig file containing the fold change in signal across the genome, which can be

**Table 2 4D Nucleome analysis pipelines.**

| Pipeline               | Steps                                   | Software                            | Available file formats | CWL/WDL filename            |
|------------------------|---|-------------------------------------|------------------------|-----------------------------|
| Hi-C <sup>1</sup>      | Alignment                               | bwa-mem                             | .bam                   | bwa-mem.cwl                 |
|                        | Filtering                               | pairtools                           | .pairs                 | hi-c-processing-bam.cwl     |
|                        | Merging replicates & matrix aggregation | cooler                              | .hic, .mcool           | hi-c-processing-pairs.cwl   |
| MARGI <sup>2</sup>     | Alignment                               | bwa-mem                             | .bam                   | imargi-processing-fastq.cwl |
|                        | Filtering                               | pairtools                           | .pairs                 | imargi-processing-bam.cwl   |
|                        | Merging replicates & matrix aggregation | cooler                              | .mcool                 | imargi-processing-pairs.cwl |
| Repli-seq <sup>3</sup> | Alignment                               | bwa-mem                             | .bam                   | repliseq-parta.cwl          |
|                        | Filtering                               | samtools                            | -                      |                             |
|                        | Binning & aggregation                   | bedtools                            | .bw, .bg               |                             |
| CUT&RUN <sup>4</sup>   | Alignment & filtering                   | bowtie2, Picard, samtools, bedtools | .bam, .bedpe           | cut-and-run-processing.cwl  |
|                        | Peak calling                            | SEACR                               | .bw, .bg, .bed         | cut-and-run-postaln.cwl     |
| ATAC-seq <sup>5</sup>  | Alignment & filtering                   | bowtie2, bedtools                   | .bed                   | atac.wdl                    |
|                        | Peak calling                            | MACS2                               | .bw, .bigbed           |                             |
| ChIP-seq <sup>6</sup>  | Alignment & filtering                   | bwa, bedtools                       | .bed                   | chip.wdl                    |
|                        | Peak calling                            | MACS2, SPP                          | .bw, .bigbed           |                             |
| RNA-seq <sup>7</sup>   | Alignment                               | STAR                                | .bam                   | rna-seq-pipeline.wdl        |
|                        | Expression quantification               | RSEM                                | .tsv                   |                             |
|                        | Read coverage                           | STAR                                | .bw                    |                             |

Listed below are (i) subdirectories for Docker images from <https://hub.docker.com/r/4ndcic>; (ii) subdirectories from github repositories at <https://github.com/4dn-dcic/> that hold the CWL or WDL pipeline files; and (iii) subdirectories for more information from <https://data.4dnucleome.org/resources/data-analysis/>. Note that for the Repliseq pipeline as well as the WDL pipelines from ENCODE, there is only one workflow file for the whole pipeline.

<sup>1</sup>4dn-hic, docker-4dn-hic/tree/v43/cwl, hi\_c-processing-pipeline.

<sup>2</sup>imargi, iMARGI-Docker/tree/v1.1.1\_dcic\_4/src/cwl, imargi-pipeline.

<sup>3</sup>repliseq, docker-4dn-repliseq/tree/v16/cwl, repli-seq-processing-pipeline.

<sup>4</sup>cut-and-run-pipeline, docker-4dn-cut-and-run-pipeline/tree/v1/cwl, cut-and-run-pipeline.

<sup>5</sup>encode-atacseq, atac-seq-pipeline, atacseq-processing-pipeline.

<sup>6</sup>encode-chipseq, chip-seq-pipeline2, chipseq-processing-pipeline.

<sup>7</sup>encode-rnaseq, rna-seq-pipeline, rnaseq-processing-pipeline information.

visualized as a 1D genome track in HiGlass, and a corresponding quality control report. The peak calling step of these processing pipelines also yields optimal peaks and conservative peaks, which are both available on the portal in bigBed format. For RNA-seq, the final output is a bigWig file that contains read counts, and can be visualized as a 1D track in HiGlass. In addition, .tsv files of gene expression and isoform expression are also available. More information is available on the portal at the following links: <https://data.4dnucleome.org/resources/data-analysis/chipseq-processing-pipeline> for ChIP-seq; <https://data.4dnucleome.org/resources/data-analysis/atacseq-processing-pipeline> for ATAC-seq; and <https://data.4dnucleome.org/resources/data-analysis/rnaseq-processing-pipeline> for RNA-seq.

Several types of quality control and assessment are performed on submitted and pipeline generated results. FastQC<sup>25</sup> is run on all fastq files, and the report generated is stored in the cloud and viewable to users. PairsQC<sup>26</sup> is a software package developed to assess the quality of .pairs files generated by the Hi-C and MARGI pipelines, and these results are also available on the portal for viewing. Additional QC reports are available on the portal for .bam alignments, as well as ChIP-seq, RNA-seq, ATAC-seq, and Repli-seq results, and report various metrics such as the percentage of reads mapped in the case of alignments, or overlap reproducibility measures in the case of ATAC-seq and ChIP-seq.

A primary advantage of using the 4DN Hi-C data is that all similar experiments have been processed in an automated fashion with the same software and versions, so the results are directly comparable and more amenable to meta-analysis. All data are processed by open source software from start (fastq files) to finish (contact matrices and domain calls). All intermediate files and full provenance graphs are available so that the user can easily find the processing steps and reproduce any portion of the full pipeline. All 4DN pipelines are available to download as Docker containers on

Docker Hub (Table 2). On the AWS cloud, data processing pipelines employ Tibanna<sup>42</sup>, developed by the 4DN DCIC. Tibanna queries metadata from the 4DN portal to obtain parameters for running the pipelines in the cloud, and updates metadata on the portal upon completion of pipelines to provide access to the processed results, which are stored in Amazon S3. For some of the other assay types, automated cloud pipelines are not yet available; in these cases, data submitters have the option to submit their own processed files that are stored and available in the cloud.

**Microscopy datasets.** Given the exciting developments in imaging technologies and their importance in investigating chromosomal dynamics, a considerable part of 4DN efforts are in the development and application of such techniques. Accordingly, the 4DN Data Portal was built to also handle microscopy datasets. Complementary to sequencing methods that measure average contact frequencies over large cell populations, imaging-based techniques such as FISH allow direct measurement of 3D distances between DNA loci in single cells. In addition to their usefulness in cross-validating findings obtained with high-throughput sequencing, microscopy data also provide valuable insight into cell-to-cell variability and dynamics that are lost in experiments with bulk samples.

The 4DN data portal hosts several types of microscopy experiments. To mention a few examples, ChromEMT combines an evolution of electron microscopy tomography with a novel DNA-labeling method, allowing visualization of chromatin organization in interphase and mitotic cells in situ at unprecedented resolution<sup>43</sup>. The OptoDroplet assay<sup>44</sup> can assess the condensation potential of proteins known to interact with membrane-less organelles in the nucleus<sup>45</sup>. Multiplexed FISH was used to obtain localization data of 8 Mb of human chromosome 19<sup>7</sup>. Single-particle tracking (SPT) experiments to



study nuclear protein dynamics<sup>46</sup> are also present on the portal. High-throughput FISH represents one example of integrating imaging and sequencing. This technique allows the examination of heterogeneity in genome organization by systematically determining the spatial position and distances between combinations of genomic interaction pairs identified by Hi-C<sup>47</sup>.

Together, the 4DN data portal hosts >600 microscopy experiment sets (Table 3) at the time of writing, divided into 27 datasets; a table is available at <https://data.4dnucleome.org/microscopy-data-overview>. Most of the data are provided as processed files: these are the most reusable and data-rich files and include, for example, locus pair distances from FISH experiments, localization files from high-throughput Single-Molecule Switching Nanoscopy<sup>48</sup>, or Spatio-temporal trajectory coordinates from SPT assays. A selection of raw image files is also available on the 4DN data portal and can be explored directly or downloaded for independent analysis. Microscopy files, unlike sequencing data,

do not undergo any automated processing or quality control measures on the data portal at this time.

**Biosamples and tiered cell lines.** Data hosted at the 4DN data portal come from a variety of biological sources. These sources are mostly human and mouse, but fly, zebrafish, chicken, hamster, and green monkey experiments can also be found. The majority of experiments were performed on cell lines, but a few are derived from tissue samples or whole organism samples. Some biosamples derive from sources such as HeLa cells, for which access needs to be restricted. For these data, we make processed results available on the data portal, but raw files are not available for download. Instead, a dbGaP identifier is linked, and the user can request access from dbGaP if necessary.



Cell lines at the 4DN data portal can be categorized into three “tiers”: Tier 1, Tier 2, and untiered. Early in the consortium, 4DN Network members had decided on a group of cell lines on which to perform coordinated experiments. The goal of this effort was to deliver data that are more directly comparable across different assays, whenever the biological question being studied would allow it. When a biosample is designated as Tier 1 in the 4DN data portal, it means that it is derived from an aliquot of cells obtained from a single common provider, to minimize sample-to-sample variation. Five human cell lines are designated as Tier 1: H1-ESC, GM12878, IMR90, HFF-hTERT (clone 6), and WTC-11. 11 different human and mouse cell lines are currently designated as Tier 2, indicating that multiple 4DN Network investigators have agreed to generate data on them. The untiered cell lines include any cell lines outside of the 4DN-designated lines. A complete list of the 4DN tiered cell lines can be found at <https://4dnucleome.org/cell-lines.html>.

**Table 3 Microscopy assay types in the 4D Nucleome Data Portal.**

| Experiment type     | No. of public experiment sets |
|---------------------|-------------------------------|
| DNA FISH            | 275                           |
| Immunofluorescence  | 138                           |
| SPT                 | 101                           |
| RNA FISH            | 77                            |
| Optodroplet         | 13                            |
| Electron Tomography | 3                             |
| Multiplexed FISH    | 1                             |
| Total               | 608                           |


The screenshot displays the 4DN data portal interface. At the top, there are buttons for 'Select All', 'Download 92 Selected Files', and 'All File Types'. Below this is a table with columns: Title, Experiment Type, Biosample, Dataset, and Condition. Two experiment sets are selected, indicated by blue checkmarks. The first is '4DNES2M5JIGV' (in situ Hi-C, H1-hESC (Tier 1)) and the second is '4DNES2R6PUEK' (in situ Hi-C, HFFc6 (Tier 1)). The second experiment set is expanded to show additional metadata, including 'Lab: Job Dekker, UMMS' and 'Publication: Krietenstein N et al. (2020) PMID:32213324'. Below the metadata, there are sections for '36 Raw Files' and '10 Processed Files'. The processed files table has columns: Experiment, File, File Type, and File Size. It lists several files, including '4DNFIJE0J221' (51.63 GB), '4DNFIFLJLISS' (17.75 GB), '4DNFIB59T7NN' (21.27 GB), '4DNFI1KESQKT' (263.11 GB), '4DNFIAIBSEAG' (27.03 GB), '4DNFIP7BS3KR' (269.29 GB), and '4DNFI7HPW10' (28.47 GB). On the left side, there is a 'Properties' sidebar with filters for 'Experiment Type' (in situ Hi-C selected) and 'Sample Tier' (Tier 1 selected).

**Fig. 1 Browsing 4DN experiment sets.** The browse view features a table of experiment sets, the second of which can be seen expanded here to show additional metadata and information about files. On the left are a number of properties that can be used to filter the results; here “in situ Hi-C” is selected as well as Tier 1 samples. The top two experiment sets in the table are also shown with their checkboxes checked so that the “Download Files” button above the table can be used.

Experiment Set | 4DNES2R6PUEK  Released | Create | Edit |  View JSON








Hi-C on HFFc6 (Tier 1) with DpnII 📅 October 24th, 2017 at 2:18pm


Source Publication >


 **Ultrastructural Details of Mammalian Chromosome Architecture.**  
Krietenstein N, Abraham S, et al., *Molecular cell* 2020

Assay Description >

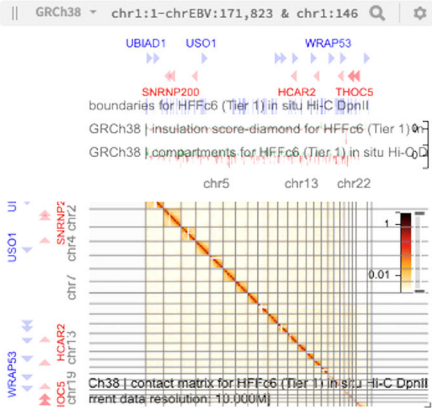
☰ Experiment Set Properties >

 10 Processed Files  36 Raw Files  Supplementary Files  Provenance  Attribution  Details  Warnings







10 Processed Files  Download 10 Processed Files

 Explore Data

GRCh38 | chr1:1-chrEBV:171,823 & chr1:146



Ch38 | contact matrix for HFFc6 (Tier 1) in situ Hi-C DpnII  
rent data resolution: 10,000M

| Experiment  | File   | File Type  | File Size                      |           |
|---|--|--|--------------------------------|-----------|
| FROM MULTIPLE EXPERIMENTS<br>4DNES2R6PUEK<br>Multiple Experiments                             | <input checked="" type="checkbox"/> 4DNFIJE0J221   | contact list-combined (pairs)                    | 51.63 GB                       |           |
|   | <input checked="" type="checkbox"/> 4DNFIFLJLISS    | contact matrix (hic)                             | 17.75 GB                       |           |
|   | <input checked="" type="checkbox"/> 4DNFI859T7NN    | contact matrix (mcool)                           | 21.27 GB                       |           |
|   | <input checked="" type="checkbox"/> 4DNFI328YQF7    | boundaries (bed)                                 | 153.59 kB                      |           |
|   | <input checked="" type="checkbox"/> 4DNFISTV2GJX    | insulation score-diamond (bw)                    | 8.04 MB                        |           |
| EXPERIMENT<br>4DNEX7P0C084<br>Bio Rep 1, Tec Rep 1<br>in situ Hi-C on HFFc6 (Tier 1) with ... | <input checked="" type="checkbox"/> 4DNFIQZ53HV     | compartments (bw)                                | 205.9 kB                       |           |
|   | <input checked="" type="checkbox"/> 4DNFI1KESQKT   | alignments (bam)                                 | 263.11 GB                      |           |
|   | <input checked="" type="checkbox"/> 4DNFIATBSEAG  | contact list-replicate (pairs)                   | 27.03 GB                       |           |
|   | EXPERIMENT<br>4DNEXRAEERUF<br>Bio Rep 2, Tec Rep 1<br>in situ Hi-C on HFFc6 (Tier 1) with ...  | <input checked="" type="checkbox"/> 4DNFIP7B5JKR | alignments (bam)               | 269.29 GB |
|   |  | <input checked="" type="checkbox"/> 4DNFI7HPMW10 | contact list-replicate (pairs) | 28.47 GB  |

**Fig. 2 Item page for a replicate set of Hi-C experiments.** A source publication is shown near the top of the page, when available. Below these are two dropdown boxes: the first, titled Assay Description, can be expanded for an explanation of the assay, and the second, titled Experiment Set Properties, contains a selection of basic metadata fields. Below is a window with several tabs. The selected tab shows the processed files associated with the experiment. The .mcool contact matrix file is visualized on the left using the integrated HiGlass browser as a 2D track, with TAD boundaries, insulation scores, and compartments as 1D tracks above it. This display can be expanded for further data exploration. Scrolling down on this page would reveal the quality control metrics associated with the processed data.

**Data Portal services.** In the following sections, we describe some of the major functionalities of the 4DN data portal. A key aspect of a useful data portal is to give an overview of the type and amount of data available and allow the user to quickly find the datasets of interest. In addition to uniformly processed raw data, extensive metadata is collected and curated, enabling searches on various aspects of the data.

**Browsing datasets.** Finding datasets can be done either by browsing or by searching. To search, one can just click the search bubble in the top navigation bar, and a bigger search box will appear. Here the user can enter search terms, and optionally choose whether they want to search all items, or search by accession or within a particular item type. To browse, the 4DN data portal homepage (<https://data.4dnucleome.org/>) has a button called “Data” in the navigation bar with the following options: Browse All, Browse Sequencing, View Microscopy, and Browse by Publication. Clicking on “Browse All” will load a page of all public Experiment Sets. On the left side of the page is a number of filtering options, where experiments can be filtered on various details. In the example shown (Fig. 1), selecting filters of “in situ Hi-C” under experiment

type and “Tier 1” under Sample Tier filters results in a smaller number. The columns labeled “Dataset” and “Condition” in the result table are useful in finding datasets of interest. Experiment Sets with the same “Dataset” name generally were generated together as part of a single study or analysis; the “Condition” field explains the differences between them. This is particularly useful when one set of replicates represents a control while another represents a treatment. Together the “dataset” and “condition” fields allow users to easily see which experiment sets go together, and what the experimental differences are.

Each result in the table can be expanded with the “+” button. The expanded view provides access to more details such as available files and the data-generating lab. Files from multiple experiment sets can be downloaded in bulk in the browse view here. Downloads require an account and a corresponding access key, but accounts can be created immediately and are free and accessible to anyone.

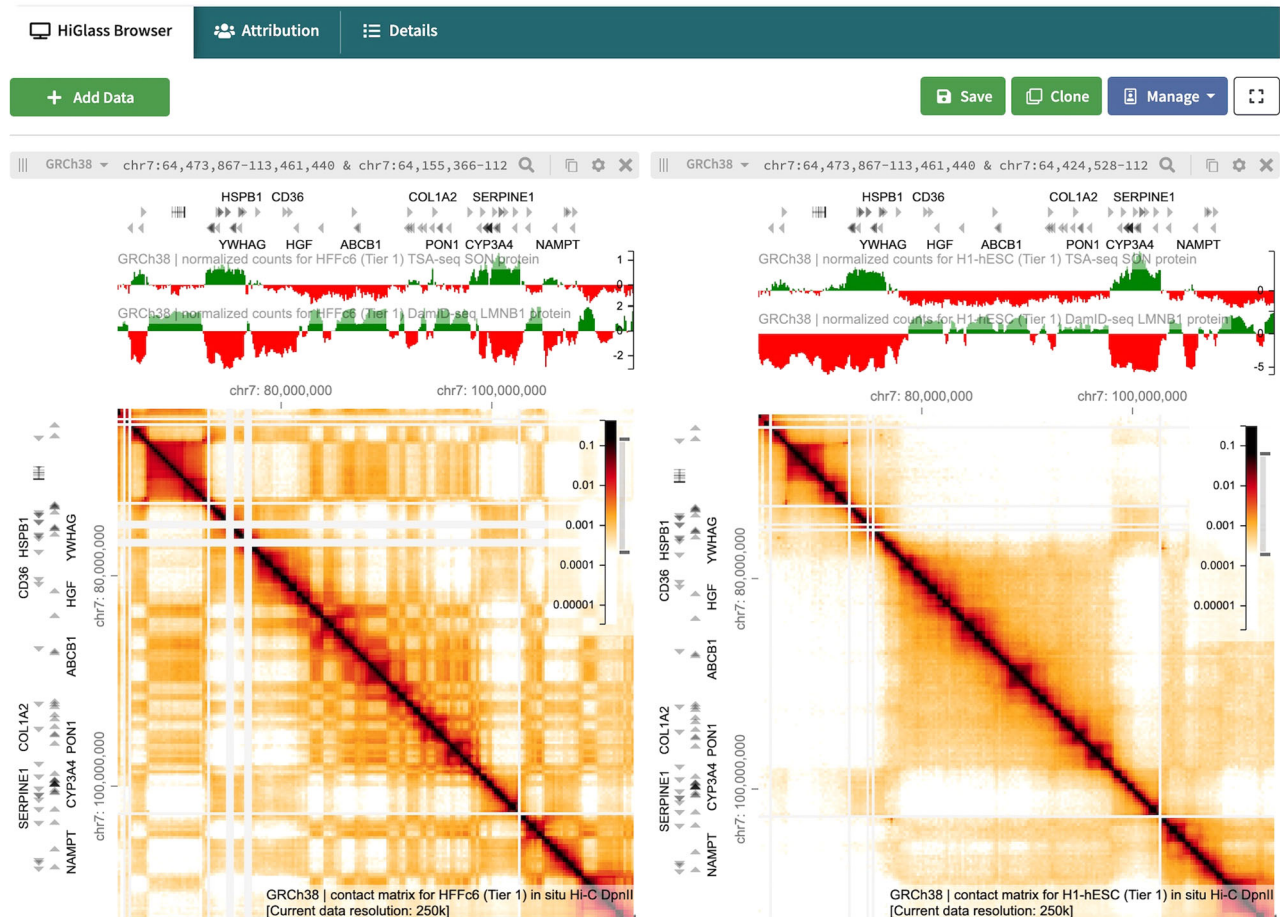
**Item pages: genomics.** Clicking on one of the accessions found in the browse results table will bring the user to that experiment set’s item page (example shown in Fig. 2). More extensive metadata is

## HiGlass Display HFF vs H1-hESC in situ Hi-C

User Content

● Draft | Edit | View JSON

📅 January 24th, 2020 at 2:40pm



**Fig. 3 HiGlass Display containing 2D and 1D Genomic Tracks.** The display shows a visualization of an HFF in situ Hi-C contact matrix on the left, and one for H1 hESCs on the right. Above each 2D contact matrix visualization are 1D tracks from a TSA-seq experiment (topmost) and a DamID-seq experiment (second from top), in the cell type that matches the corresponding matrix. The “Add Data” button in the top left gives the user the option to add more files to the visualization display; on the top right, the user has the ability to save the display, clone the display to create a new item, or to manage permissions of who can view the display.

available on these pages. A reference publication will be featured at the top if the dataset has been published. If the dataset has not been published, it is still available for use, with the following data usage guidelines: (1) the data-generating lab should be contacted to discuss possible coordinated publication; (2) the 4DN white paper<sup>1</sup> should be cited; and (3) the lab which generated the data should be acknowledged. Further below on the page is a pane with several tabs, each providing more details about a different aspect of the experiment set. Typically these include Processed Files, Raw Files, Provenance, Attributions, Details, and Warnings or Commendations. The Raw Files and Processed Files tabs provide metadata about files associated with the experiment set and also designate the replicate from which they were generated. For processed files, this tab often includes a small HiGlass view of the final processed output, which can be expanded into a new window. If processed files were generated by one of the 4DN processing pipelines, there will also be a provenance graph available to view in a subsequent tab. This directed graph shows the inputs and outputs for each step in the pipeline that was run on the experiment set. In the case of user-submitted processed files where a 4DN processing pipeline

was not run, a provenance graph will still be present, but there will only be a single step titled “File Provenance Tracking Workflow” and it will show only relationships between input and output files, as the portal does not have any metadata about details of the processing performed by the user.


The final tab on item pages for Experiment Sets, Experiments, and Biosamples will be Warnings or Commendations. Presently these are used to indicate if an experiment lacks replicates, or to indicate whether the biosamples have all of the requisite metadata. The biosamples with a gold commendation have extensive metadata information as well as a morphology image taken before harvesting. This allows people who want to reuse the data to know more about the cultures and whether they had been growing with typical morphology. If a sample does not meet one or more metadata requirements for the gold commendation, this will show up under Warnings with a message indicating which piece(s) of metadata is missing.

All item types will also have Attributions and Details tabs. The attribution tab indicates which lab(s) generated the data and any associated publications. Datasets published in journals often must


## Replicate Experiments

Experiment Set | 4DNESA84SNKC  Released | [Create](#) | [Edit](#) | [View JSON](#)

DNA-FISH with probe pPK1003

 December 2nd, 2019 at 12:07pm

Source Publication &gt;

 [Distinct features of nucleolus-associated domains in mouse embryonic stem cells](#)  
Aizhan Bizhanova, Aimin Yan, et al., *bioRxiv* 2019

Assay Description &gt;

 Experiment Set Properties Experiment Set Type 

Replicate

Organism

[M. musculus](#)

Biosource Type

stem cell derived cell line

Biosource

F121-9-CASTx129 (Tier 2)

Experiment Type(s)

[DNA FISH](#)

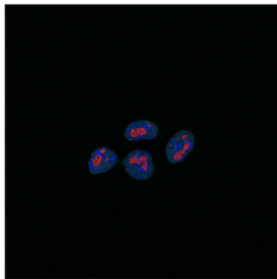
Modification Type

None

Treatment Type

None

Assay Details

Target: Nucleolus, Chromosomes,  
GRCm38:15:84216729-84409126 mouse  
regionSample Image Sample Image - DNA FISH (green) with  
Fibrillarlin (nucleolus, red) and DAPI  
(blue) co-staining[View File Item - 4DNFIOFUTRNO](#)

Imaging Paths

|      |  |
|------|--|
| ch00 | <a href="#">Nucleolus targeted by Alexa A594-labeled Anti-Rabbit Secondary Antibody (with Rabbit Anti-Fibrillarlin Antibody)</a> |
| ch01 | <a href="#">GRCm38:15:84216729-84409126 mouse region targeted by Alexa 488-labeled Streptavidin (with Biotin)</a>                |
| ch02 | <a href="#">Chromosomes targeted by DAPI</a>   |

 3 Processed Files 15 Raw Files Provenance Attribution Details Warnings

**Fig. 4 Microscopy experiment set item page.** Item pages for microscopy experiment sets are similar to those for genomics, but with a few important differences. The sample image field contains an image preview, which can be clicked for a popup containing an interactive image viewer. Beside the sample image is a field called “Imaging Paths”, which describes what is imaged in each channel for the experiments, including the biological target and the antibodies or probes used.

also be deposited in GEO or SRA, and some 4DN datasets have also been deposited to the ENCODE portal; external datasets were often originally deposited at GEO and SRA before being added to 4DN. If available, identifiers to other databases that may also host the data are also included here. The details tab has a list of all metadata fields associated with the item, regardless of whether or not it is already displayed elsewhere on the page.

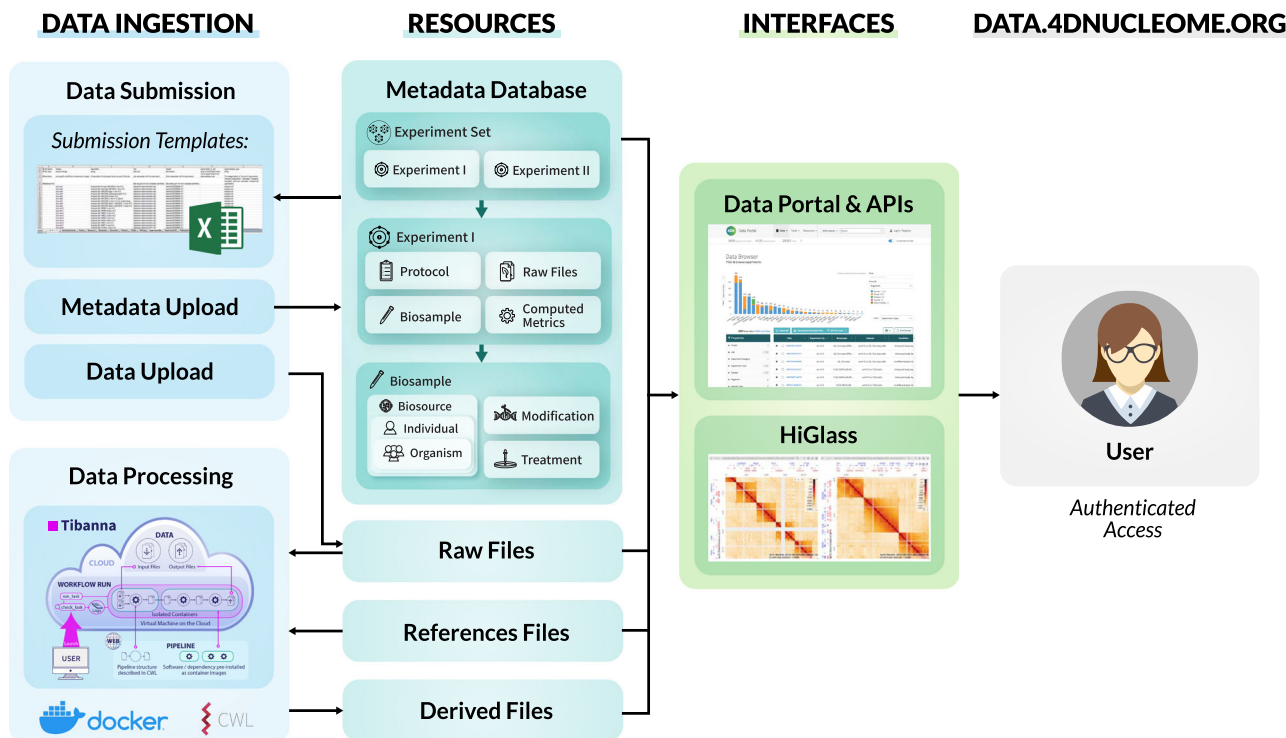
**Comparing files in HiGlass.** Whereas several genomic and epigenomic data repositories have a 1D genome browser tool, contact frequencies across distant genomic locations require 2D visualization. To this end, a HiGlass browser<sup>8</sup> has been integrated into the portal for interactive visualization of both 2D contact matrices and 1D genomic tracks. The 4DN visualization workspace can be accessed either at <https://data.4dnucleome.org/tools/visualization> or by clicking on the “Explore Data” button above the miniature HiGlass view displayed on an Experiment Set item

page. A short tutorial video can also be found at <https://www.youtube.com/watch?v=LEDaOa3NZtM>.

The close integration of HiGlass into the portal allows users to leverage the portal search capabilities to identify files of interest for visualization that can then be added to existing HiGlass views. Files can be added to the HiGlass display by clicking on the Add Data button, which opens a window where users can filter via various metadata fields to find a file of interest. Multiple 2D contact matrices can be compared at once (Fig. 3), and the position and zoom level of each view are locked to the first one, such that zooming in or out or dragging the mouse to different genomic locations moves all the matrices in unison.

1D tracks can also be added to each HiGlass view. By default, each 1D track added will be added to every view in the display. This allows it to be compared to each 2D matrix simultaneously. Multiple 1D tracks can be added in this fashion. Users can then save their customized views for later retrieval and sharing with collaborators.





**Fig. 5 4DN Data Portal architecture.** Metadata is submitted by users via spreadsheet forms and gets loaded into the database, while associated files are uploaded into cloud storage. Once metadata and files are in place, automated processing pipelines can be run on AWS using the Tibanna pipeline runner. All metadata can then be searched via the data portal website, where files can also be visualized with HiGlass. The data portal website is accessible to external users, who can also login in order to download files.

**Finding microscopy experiments.** In the Data menu in the top navigation bar, clicking on “View Microscopy Datasets” will bring the user to a table of the microscopy datasets that can be found in the portal. These are organized by Dataset label rather than by Experiment Set, as some datasets comprise large numbers of replicate sets, allowing the table to be much more concise and straightforward. The second column contains a link to a browse page of all the Experiment Sets in the dataset.

A microscopy experiment set item page (Fig. 4) is organized much the same as the sequencing experiment sets, with a few notable differences. In the overview section above the tabs, in addition to common metadata, there may be two extra fields: a sample image, for sets that include raw files, and information about the imaging paths. The Imaging Path refers to metadata about what was imaged, including cellular target, antibodies or probes, and the detected label. In addition, there are no processing pipelines run by the 4DN data portal on microscopy datasets, so the Provenance tab will often be dimmed indicating no provenance is available. However, in some cases submitting laboratories indicate which files are processed results derived from others and in these cases, a simple provenance graph reflecting this is shown.

If a sample image is present and derives from a raw image file on the portal rather than a rendered jpg, then clicking on the image will bring the user to an interactive image viewer, where it is possible to scroll through focal planes in a z-stack, or through time points in a time-lapse image, or adjust signal and color levels in the display.

**Discussion**

With 1273 sequencing and 608 microscopy experiment sets that are publicly available and consistently curated, the portal provides a high volume of data in a wide array of data types to nuclear

architecture researchers across the globe. The portal interface provides user-friendly search capabilities with many options for filtering results, and the result pages have been designed to communicate extensive metadata for each item. Uniform processing on assays with automated pipelines also ensures that the processed results are consistent and directly comparable.

The integrated HiGlass browser allows the data portal to serve not only as a hub for searching and downloading research data but also as a visualization platform. 2D output from genomic assays like Hi-C can be directly compared to other 2D assays like DNA SPRITE, or to results from 1D assays like ChIP-seq or CUT&RUN. Even when only looking at Hi-C experiments, some 4DN data contributors have uploaded experiments comparing different variations in Hi-C protocols<sup>11</sup>, and the HiGlass browser provides an easy way to examine differences in data output resulting from protocol modifications.

Prior to the development of the 4DN data portal, there was a lack of a centralized repository that specializes in chromatin conformation and nuclear structure data, specifically. Other repositories that host Hi-C and similar datasets include the ENCODE portal<sup>49</sup>, as well as NCBI’s GEO<sup>50</sup> and EMBL-EBI’s ArrayExpress<sup>51</sup>. GEO and ArrayExpress are great resources for hosting published sequencing datasets, but they lack a finer specialized focus or visualization capabilities for processed results. Their metadata also tends to be minimal and less structured, which is understandable given the huge diversity of datasets that need to be hosted there. ENCODE has an infrastructure that is very similar to the 4DN data portal, and it has amassed an impressive number of sequencing datasets over the years it has been active. Although they have Hi-C datasets available, since their portal has been engineered to focus more on 1D datasets, they lack a 2D genome browser and the Hi-C datasets cannot be visualized there. The 4DN data portal thus fills a very important niche: it represents a centralized repository for nuclear structure

data specifically, with a data model tailored to provide extensive metadata about each experiment; it enables visualization of processed results, where 2D genomic data can be compared to 1D genomic tracks; and it also houses microscopy data that can complement sequencing results. To our knowledge, this is unique to the 4DN data portal.

The second phase of the NIH Common Fund 4D Nucleome Project began in 2020, and will proceed for five years, with a focus on the role of nuclear structure and function in human health and disease, as well as on the continued development of data visualization and integration tools. Thus, we expect the data volume in the 4DN data portal to continue to increase. Development on the portal is ongoing to ensure that the 4DN data portal continues to serve the needs of the nuclear biology research community.

## Methods

**Portal architecture.** The 4DN Consortium was projected to generate large volumes of diverse datasets at its initiation. To be useful to scientists with different technical and scientific backgrounds, the primary objectives were to ensure users can determine what data are available on the portal, and that data of different modalities should be made accessible to all scientists. This requires an architecture that is modular, responsive, and scalable, with an easily extensible data model. Our software architecture is entirely cloud-based, as we envisioned that taking analytical tools to the cloud environment where data reside, rather than downloading large amounts of data to a local server, is quickly becoming the preferred mode of data analysis.

The design of the 4DN data portal infrastructure (Fig. 5), originally based on the ENCODE infrastructure, includes the following components: (1) A postgres database storing metadata in json format, first developed in ENCODE; (2) The python pyramid framework for the database known as SnoVault<sup>52</sup>, first developed in ENCODE but further tailored and developed by the 4DN DCIC (<https://github.com/4dn-dcic/snovault>); (3) The FourFront front-end (<https://github.com/4dn-dcic/fourfront>), originally based on EncodeD<sup>52</sup> from ENCODE, but engineered by 4DN DCIC to feature a data model for representing diverse datasets, and includes a modern front-end with reactJS to provide a responsive user experience; (4) Elasticsearch that provides fast and efficient search with various metadata fields by indexing all items and formatting them for retrieval; (5) AWS S3 used for file storage, enabling all public data files to be accessed via the data portal interface; (6) A RESTful API underlying the infrastructure, through which all metadata in the portal can be accessed.

In the AWS cloud environment, changes in the database are indexed within seconds to minutes, allowing updates and releases of datasets with minimal overhead. Similarly, the infrastructure was developed with features that allow the release of changes to a data model and software versions overnight without any server downtime. Thus, we are able to ingest and release datasets with new data models at any time, without having to enforce data release cycles.

All of the software used in the data portal is open access and open source (<https://github.com/4dn-dcic/>). Finally, the 4DN data portal also has an associated support contact email ([support@4dnucleome.org](mailto:support@4dnucleome.org)), allowing curators to quickly address any question or concern from end-users.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

All datasets described are available at <https://data.4dnucleome.org/>.

## Code availability

All of the software that comprises the 4DN data portal infrastructure is free and open source. All code repositories mentioned are available from <https://github.com/4dn-dcic>.

Received: 17 October 2021; Accepted: 28 March 2022;

Published online: 02 May 2022

## References

- Dekker, J. et al. The 4d nucleome project. *Nature* **549**, 219–226 (2017).
- Lieberman-Aiden, E. et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science* **326**, 289–293 (2009).
- Hansen, R. S. et al. Sequencing newly replicated DNA reveals widespread plasticity in human replication timing. *Proc. Natl. Acad. Sci. USA* **107**, 139–144 (2010).
- van Steensel, B. & Henikoff, S. Identification of in vivo DNA targets of chromatin proteins using tethered dam methyltransferase. *Nat. Biotechnol.* **18**, 424–428 (2000).
- Skene, P. J. & Henikoff, S. An efficient targeted nuclease strategy for high-resolution mapping of DNA binding sites. *eLife* **6**, e21856 (2017).
- Quinodoz, S. A. et al. Higher-order inter-chromosomal hubs shape 3d genome organization in the nucleus. *Cell* **174**, 744–757.e24 (2018).
- Nir, G. et al. Walking along chromosomes with super-resolution imaging, contact maps, and integrative modeling. *PLoS Genet.* **14**, e1007872 (2018).
- Kerpedjiev, P. et al. HiGlass: web-based visual exploration and analysis of genome interaction maps. *Genome Biol.* **19**, 125 (2018).
- Hsieh, T.-H. S. et al. Mapping nucleosome resolution chromosome folding in yeast by micro-c. *Cell* **162**, 108–119 (2015).
- Deng, X. et al. Bipartite structure of the inactive mouse x chromosome. *Genome Biol.* **16**, 152 (2015).
- Akgol Oksuz, B. et al. Systematic evaluation of chromosome conformation capture assays. *Nat. Methods* **18**, 1046–1055 (2021).
- Nagano, T. et al. Single-cell hi-c reveals cell-to-cell variability in chromosome structure. *Nature* **502**, 59–64 (2013).
- Ramani, V. et al. Massively multiplex single-cell hi-c. *Nat. Methods* **14**, 263–266 (2017).
- Flyamer, I. M. et al. Single-nucleus hi-c reveals unique chromatin reorganization at oocyte-to-zygote transition. *Nature* **544**, 110–114 (2017).
- Fullwood, M. J. et al. An oestrogen receptor alpha-bound human chromatin interactome. *Nature* **462**, 58–64 (2009).
- Li, X. et al. Long-read ChIA-PET for base-pair resolution mapping of haplotype-specific chromatin interactions. *Nat. Protoc.* **12**, 899–915 (2017).
- Zheng, M. et al. Multiplex chromatin interactions with single-molecule precision. *Nature* **566**, 558–562 (2019).
- Fang, R. et al. Mapping of long-range chromatin interactions by proximity ligation-assisted ChIP-seq. *Cell Res.* **26**, 1345–1348 (2016).
- Sridhar, B. et al. Systematic mapping of RNA-chromatin interactions in vivo. *Curr. Biol.* **27**, 602–609 (2017).
- Yan, Z. et al. Genome-wide colocalization of RNA-DNA interactions and fusion RNA pairs. *Proc. Natl. Acad. Sci. USA* **116**, 3328–3337 (2019).
- Wu, W. et al. Mapping RNA-chromatin interactions by sequencing with iMARGI. *Nat. Protoc.* **14**, 3243–3272 (2019).
- Chen, Y. et al. Mapping 3d genome organization relative to nuclear compartments using TSA-seq as a cytological ruler. *J. Cell Biol.* **217**, 4025–4048 (2018).
- Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint at <https://arxiv.org/abs/1303.3997> (2013).
- Goloborodko, A., Abdennur, N., Venev, S., Brandao, H. & Fudenberg, G.. Pairtools v0.3.0. *zenodo* <https://doi.org/10.5281/zenodo.2649383> (2019).
- Andrews, S. FastQC: a quality control tool for high throughput sequence data. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>. (2010).
- Lee, S., Vitzthum, C., Alver, B. H. & Park, P. J. Pairs and pairix: a file format and a tool for efficient storage and retrieval for hi-c read pairs. *Bioinformatics* **37**, btab870 (2021).
- Abdennur, N. & Mirny, L. A. Cooler: scalable storage for hi-c data and other genomically labeled arrays. *Bioinformatics* **36**, 311–316 (2020).
- Durand, N. C. et al. Juicebox provides a visualization system for hi-c contact maps with unlimited zoom. *Cell Syst.* **3**, 99–101 (2016).
- Dixon, J. R. et al. Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature* **485**, 376–380 (2012).
- Nora, E. P. et al. Spatial partitioning of the regulatory landscape of the x-inactivation center. *Nature* **485**, 381–385 (2012).
- Beagan, J. A. & Phillips-Cremins, J. E. On the existence and functionality of topologically associating domains. *Nat. Genet.* **52**, 8–16 (2020).
- Szabo, Q., Bantignies, F. & Cavalli, G. Principles of genome folding into topologically associating domains. *Sci. Adv.* **5**, eaaw1668 (2019).
- Zufferey, M., Tavernari, D., Oricchio, E. & Ciriello, G. Comparison of computational methods for the identification of topologically associating domains. *Genome Biol.* **19**, 217 (2018).
- Venev, S. et al. Cooltools: v0.4.1. *zenodo* <https://doi.org/10.5281/zenodo.5214125> (2021).
- Crane, E. et al. Condensin-driven remodelling of x chromosome topology during dosage compensation. *Nature* **523**, 240–244 (2015).
- Li, H. et al. The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
- Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
- Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
- The Broad Institute. *Picard toolkit*. <https://broadinstitute.github.io/picard/> (2019).

40. Meers, M. P., Tenenbaum, D. & Henikoff, S. Peak calling by sparse enrichment analysis for CUT&RUN chromatin profiling. *Epigenetics Chromatin* **12**, 42 (2019).
41. Davis, C. A. et al. The encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic Acids Res.* **46**, D794–D801 (2018).
42. Lee, S. et al. Tibanna: software for scalable execution of portable pipelines on the cloud. *Bioinformatics* **35**, 4424–4426 (2019).
43. Ou, H. D. et al. ChromEMT: visualizing 3d chromatin structure and compaction in interphase and mitotic cells. *Science* **357**, eaag0025 (2017).
44. Shin, Y. et al. Spatiotemporal control of intracellular phase transitions using light-activated optoDroplets. *Cell* **168**, 159–171.e14 (2017).
45. Courchaine, E. M. et al. DMA-tudor interaction modules control the specificity of in vivo condensates. *Cell* **184**, 3612–3625.e17 (2021).
46. Hansen, A. S. et al. Robust model-based analysis of single-particle tracking experiments with spot-on. *Elife* **7**, e33125 (2018).
47. Finn, E. H. et al. Extensive heterogeneity and intrinsic variation in spatial genome organization. *Cell* **176**, 1502–1515.e10 (2019).
48. Barentine, A. E. S. et al. An integrated platform for high-throughput nanoscopy. Preprint at <https://www.biorxiv.org/content/10.1101/606954v3> (2019).
49. Sloan, C. A. et al. ENCODE data at the ENCODE portal. *Nucleic Acids Res.* **44**, D726–D732 (2016).
50. Barrett, T. et al. NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res.* **41**, D991–995 (2013).
51. Parkinson, H. et al. ArrayExpress update—an archive of microarray and high-throughput sequencing-based functional genomics experiments. *Nucleic Acids Res.* **39**, D1002–1004 (2011).
52. Hitz, B. C. et al. SnoVault and encodeD: A novel object-based storage system and applications to ENCODE metadata. *PLoS One* **12**, e0175310 (2017).
53. Mifsud, B. et al. Mapping long-range promoter contacts in human cells with high-resolution capture hi-c. *Nat. Genet.* **47**, 598–606 (2015).
54. Kaya-Okur, H. S. et al. CUT&Tag for efficient epigenomic profiling of small samples and single cells. *Nat. Commun.* **10**, 1390 (2019).
55. Marchal, C. et al. Genome-wide analysis of replication timing by next-generation sequencing with *e1* repli-seq. *Nat. Protoc.* **13**, 819–839 (2018).
56. Akhtar, W. et al. Chromatin position effects assayed by thousands of reporters integrated in parallel. *Cell* **154**, 914–927 (2013).
57. Vertii, A. et al. Two contrasting classes of nucleolus-associated domains in mouse fibroblast heterochromatin. *Genome Res.* **29**, 1235–1249 (2019).
58. Beagrie, R. A. et al. Complex multi-enhancer contacts captured by genome architecture mapping. *Nature* **543**, 519–524 (2017).
59. Belaghzal, H. et al. Liquid chromatin hi-c characterizes compartment-dependent chromatin interaction dynamics. *Nat. Genet.* **53**, 367–378 (2021).
60. Paulsen, M. T. et al. Use of bru-seq and BruChase-seq for genome-wide assessment of the synthesis and stability of RNA. *Methods* **67**, 45–54 (2014).
61. Tavares-Cadete, F., Norouzi, D., Dekker, B., Liu, Y. & Dekker, J. Multi-contact 3c reveals that the human genome during interphase is largely not entangled. *Nat. Struct. Mol. Biol.* **27**, 1105–1114 (2020).

## Acknowledgements

We would like to thank the current and past members of the 4DN Network for all the time and effort they put into data generation, and without whom this data portal would not have been built. We also thank J. Michael Cherry and the development team at ENCODE for helpful discussions and advice regarding setting up and running a data portal, as well as for allowing us to build off of the ENCODE software ecosystem. This work was funded by the NIH Common Fund grant 1U01CA200059 to P.J.P.

## Author contributions

S.B.R. wrote the manuscript with the help of A.J.S., A.C., K.K., C.B., L.M., B.H.A., and P.J.P. J.J., A.K.B., C.V., S.B.R., A.J.S., K.K., A.C., C.B., L.M., S.L., A.D.V., W.R., K.M.P., P.K., N.A., M.L., S.U.O., and U.C. developed the software. A.J.S., K.K., A.C., S.B.R., and L.M. handled data submissions. S.R.E. carried out graphic design. P.J.P. supervised the work along with N.G., L.A.M., and B.H.A.

## Competing interests

The authors declare no conflict of interest.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-022-29697-4>.

**Correspondence** and requests for materials should be addressed to Peter J. Park.

**Peer review information** *Nature Communications* thanks the anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022, corrected publication 2022