

RESEARCH ARTICLE

# RNAprecis: Prediction of full-detail RNA conformation from the experimentally best-observed sparse parameters

Henrik Wiechers<sup>1</sup>\*, Christopher J. Williams<sup>2</sup>, Benjamin Eltzner<sup>1</sup>, Franziska Hoppe<sup>1</sup>, Michael G. Prisant<sup>2</sup>, Vincent B. Chen<sup>2</sup>, Ezra Miller<sup>3</sup>, Kanti V. Mardia<sup>4</sup>, Jane S. Richardson<sup>2\*</sup>, Stephan F. Huckemann<sup>1\*</sup>

**1** Felix-Bernstein-Institute for Mathematical Statistics in the Biosciences, University of Göttingen, Göttingen, Germany, **2** Department of Biochemistry, Duke University School of Medicine, Durham, North Carolina, United States of America, **3** Mathematics Department, Duke University, Durham, North Carolina, United States of America, **4** Department of Statistics, University of Leeds, Leeds, United Kingdom

\* These authors contributed equally to this work.

\* [jsr@kinemage.biochem.duke.edu](mailto:jsr@kinemage.biochem.duke.edu) (JSR); [Stephan.Huckemann@mathematik.uni-goettingen.de](mailto:Stephan.Huckemann@mathematik.uni-goettingen.de) (SFH)



**OPEN ACCESS**

**Citation:** Wiechers H, Williams CJ, Eltzner B, Hoppe F, Prisant MG, Chen VB, et al. (2026) RNAprecis: Prediction of full-detail RNA conformation from the experimentally best-observed sparse parameters. *PLoS Comput Biol* 22(5): e1014242. <https://doi.org/10.1371/journal.pcbi.1014242>

**Editor:** Peter M Kasson, Georgia Institute of Technology, UNITED STATES OF AMERICA

**Received:** October 23, 2025

**Accepted:** April 15, 2026

**Published:** May 8, 2026

**Copyright:** © 2026 Wiechers et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data availability statement:** The code used to generate the analyses and plots presented in this paper can be found <https://gitlab.gwdg.de/henrik.wiechers1/rnaprecis> Data used is available at <https://doi.org/10.5281/zenodo.8103014>.

## Abstract

We address the problem of predicting high-detail RNA structure geometry from the information available in low-detail experimental maps. Here, low-detail refers to resolutions  $\approx 2.5\text{-}3.5\text{\AA}$ , where the location of the phosphate groups and the glycosidic bonds can be determined from experimental maps but all other backbone atom positions cannot. In contrast, higher-resolution maps allow high-detail determinations of all backbone atomic positions. To this end, we first create a gold standard dataset of highly curated, experimentally supported RNA suites. Second, we develop and employ a modified version of the previously devised algorithm MINT-AGE to learn clusters that are in high correspondence with the gold standard's conformational classes of suites based on 3D RNA structure. Since some of the gold standard classes are of very small size, a new modified version of MINT-AGE is able to also identify very small clusters. Third, we create a new conformer prediction algorithm, RNAprecis, which assigns low-detail structures to newly designed 3D shape coordinates. Our improvements include: (i) learned classes augmented to cover also very low sample sizes and (ii) replacing distances from clusters by Bayesian posterior probabilities. On test data containing suites modeled as conformational outliers, RNAprecis shows good results suggesting that our learning method generalizes well. In particular, we show that the modified MINT-AGE clustering can more finely delineate between previously unseen suite conformer separations. For example, the **0a** conformer has been separated into two clusters seen in different structural contexts. Such new distinctions can have implications for biochemical interpretation of RNA structure.

**Funding:** BE, SH and HW acknowledge support from Deutsche Forschungsgemeinschaft SFB 1456; <https://www.dfg.de> SH acknowledges support from Deutsche Forschungsgemeinschaft HU-1575-7, Deutsche Forschungsgemeinschaft GK 2088 and the Niedersachsen Vorab of the Volkswagen Foundation; <https://www.volkswagenstiftung.de> CJW, MGP, VBC, and JSR acknowledge support from National Institutes of Health P01-GM063210; <https://www.nih.gov> and the Phenix Industrial Consortium; <https://phenix-online.org/consortium> EM acknowledges support from U. S. National Science Foundation grant DMS-2515765; <https://www.nsf.gov> The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

## Author summary

This article addresses the problem that RNA molecular structures inferred from measurement and stored in molecular structure data bases frequently contain errors on the atomic scale, leading to chemically impossible molecule shapes. We provide a method to identify valid RNA structure elements from chemically viable structures and use a statistical comparison on a coarse grained scale to indicate correct atom placement. On a data set of faulty structures with known corrections, we achieve a high rate of matching corrections, showing the power of our method.

## 1 Introduction

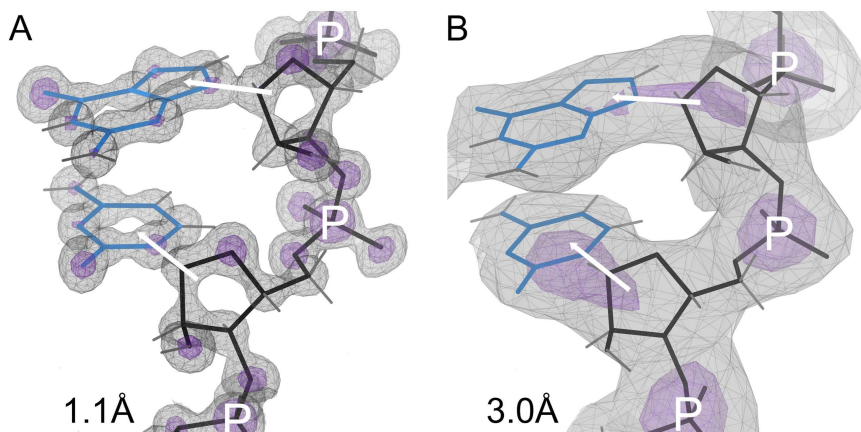
In this work, we develop a new system, RNAprecis, for analysis of RNA backbone structure. To do this, we produce a quality-curated training set for automated use in analyzing the increasingly available number of RNA structures and develop a modified version of MINT-AGE. As a multidisciplinary collaboration between structural biochemistry and mathematical statistics, we adopt sophisticated treatments from non-Euclidean statistics, and we will cover a range of background material below.

### 1.1 Prior work and background in structural biology

A large fraction of RNA molecules form well ordered tertiary structures that have specific binding and catalytic functions. Achieving an accurate model of the RNA backbone is important for understanding RNA structure and function, because in RNA many of the significant functional contacts are with the backbone, as in the peptide transfer reaction of protein synthesis in ribosomes [1].

Correctly modeling RNA backbone conformations in X-ray crystallography and cryoEM experimental 3D structures is a long-standing challenge in structural biology. The RNA backbone has many degrees of freedom: 7 dihedral angles, plus the ribose sugar pucker. At “atomic” resolutions better than about 1.5Å, the positions of all individual atoms in the atomic model are experimentally determined, as by the electron density map depicted in Fig 1A. At resolutions worse than 2-2.5Å, the map for the backbone gradually becomes more featureless and ambiguous. The model’s fit to the map becomes less able to distinguish correct from incorrect local conformations, compounding the challenge presented by RNA’s many degrees of freedom.

As shown in Fig 1B, at resolutions of 3-3.5Å, the phosphate positions and glycosidic bond directions can still be observed directly (and they may be the only features that can be reliably observed). The phosphorus is a much heavier atom than the rest, with a strong map signal, and is well located across this range of resolutions because it is the center of the symmetric PO<sub>4</sub> group (the purple contours in Fig 1). The nucleobase is large, planar, and often held in place by base-pairing, and it usually produces a recognizable blob in the map, while the sugar forms a 3-way blob where the base and backbone join.



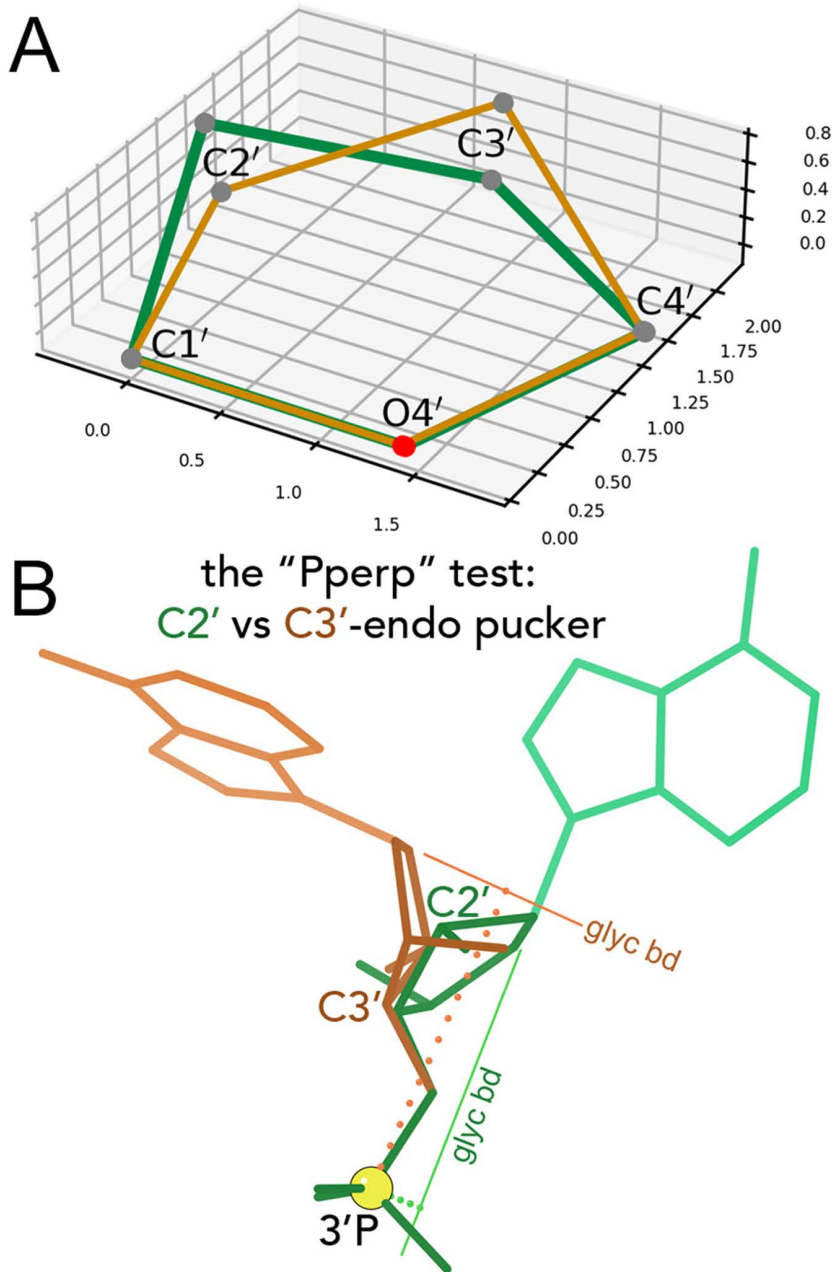
**Fig 1. The positions of individual atoms in the derived atomic model are experimentally determined by the electron density when at high resolution (Panel A; 1.1Å), but not at low resolution (Panel B; 3Å).** However, the phosphate positions and glycosidic bond directions can be seen directly at low resolution, and the pucker can be determined from them. Electron density contours are at  $1.2\sigma$  (gray) and  $3\sigma$  (purple) where  $\sigma$  is the standard deviation of the corresponding electron density map. Models have backbone in black, bases in blue. Phosphates are marked with P, glycosidic bond directions with a white arrow between sugar and base.

<https://doi.org/10.1371/journal.pcbi.1014242.g001>

In contrast, the pucker state of the 5-membered ribose ring is hidden inside a large featureless blob at resolutions of 3–3.5Å. In RNA the pucker is quite cleanly bimodal (Fig 2A), as either C3'-endo (most common, and in A-form) or C2'-endo (frequent in irregular or non-helical regions). Other puckers are extremely rare in RNA. (DNA, in contrast, has many deoxyribose puckers.) The ribose pucker state is a dominant influence on RNA conformation, because a pucker change greatly alters the relative directions of the three large ring substituents: the base, the preceding backbone, and the following backbone (Fig 1B).

Almost 20 years ago [2], it was discovered that the ribose sugar pucker can be successfully and simply deduced from the more reliable observations of phosphorus position and glycosidic bond direction, even when the pucker itself has been modeled incorrectly. As shown in Fig 2B, the length of a perpendicular dropped from the next-in-sequence (3') phosphorus atom to the extended vector direction of the glycosidic bond that connects base to ribose is  $< 2.9\text{\AA}$  for C2'-endo and  $\geq 2.9\text{\AA}$  for C3'-endo. This process for identifying pucker was incorporated into MolProbity validation [2,3], and the measure was named as the “Pperp” (Phosphate *per*pendicular) criterion [4]. Correcting a sugar pucker based on this deduction consistently corrects clashes and bond length and angle outliers within its suite, removes most difference density peaks, and in the aggregate lowers the main residual R and the residual  $R_{\text{free}}$  calculated with omitted data [4–6]. It works well with more complex systems such as ERRASER [7,8], and it has been successfully in use within Phenix refinement as well as validation for over 15 years [9,10], during which time Phenix has been disproportionately often used for RNA structures. There are other errors that cause outliers in RNA, but sugar pucker is a major one.

These map features and challenges generally hold across both X-ray crystallography's electron density maps and single-particle cryo Electron Microscopy's (cryoEM's) Coulomb potential maps, the two predominant experimental methods for solving 3D structures of RNA. Sample preparation, data collection, and data processing are completely different, but both methods produce quite similar-looking 3D maps. CryoEM is sensitive to charge as well as mass, so for nucleic acids the positive bases are somewhat stronger and the negative phosphates somewhat weaker, but not to a degree that affects the Pperp criterion. In most experimental RNA structures, the overall resolution applies to a large majority of the residues or suites; however, local or effective resolution can vary. As local mobility or uncertainty increases, the map becomes less distinct, then patchier, and eventually disappears into noise. In such regions, the assumption of fairly robust  $\text{PO}_4$  position and glycosidic bond direction becomes inapplicable.



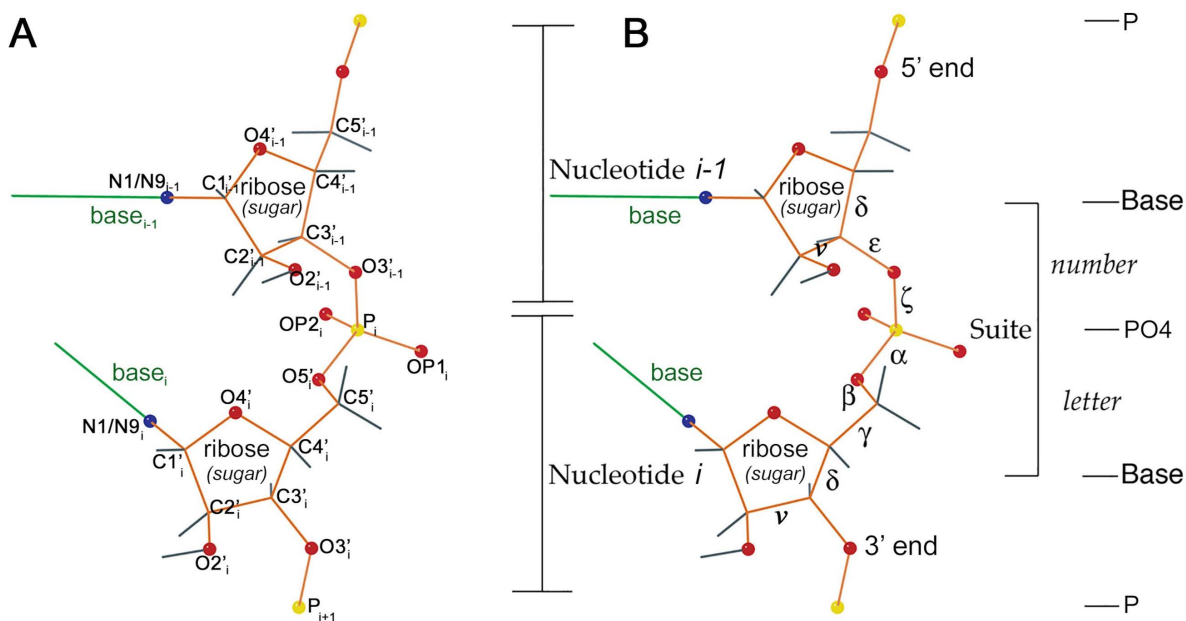
**Fig 2. Illustration of sugar puckers and the Pperp criterion.** A: The two RNA ribose sugar puckers. In a C3' endo pucker (brown), the C3' atom is above (in positive orientation) the oriented plane spanned by C4'-O4' and C1'-O4'. In a C2' endo pucker (green), the C2' atom is above the plane. Typically, all atoms but the puckered one are (nearly) co-planar. B: Contrasted cases, superimposed at their phosphate ends, to illustrate the Pperp criterion for distinguishing between C3'-endo (brown) and C2'-endo (green) sugar puckers in RNA. The atom puckered out of plane in each ring is labeled, and the vector extension of each glycosidic bond is drawn and labeled. The next P atom in sequence is marked by a large yellow ball. Dotted lines show the perpendiculars dropped from the P to each of the glycosidic bond vectors: long (and brown) for the C3'-endo puckered case; short (and green) for the C2'-endo case.

<https://doi.org/10.1371/journal.pcbi.1014242.g002>

The overall success of the Pperp pucker analysis made us hope that other combinations of parameters, taken only from the best-observed features, might enable deduction of full-detail RNA backbone conformation from low-detail representations. We chose the name RNAprecis for this desired system, to indicate a concise description that captures the important aspects. (Precis: *pray'*-see, in writing composition, a short formal summary of a much longer document.)

**Fig 3** diagrams the main features of full-detail local RNA 3D structure that we describe and analyze in this work. In Panel A the atomic model, with atoms named, covers two phosphate-to-phosphate nucleotide residues and one central sugar-to-sugar segment, known as a *suite*. For describing RNA backbone conformations, this suite is a more useful division than the chemical residue from P to P. The suite relates adjacent bases, including such important interactions as base stacking, and its atoms have many more contacts with one another than inside a residue [11]. A suite includes atoms from two adjacent, bonded residues. By convention, the suite is named by the residue number (*i*) of its second member residue—the same residue number as the P atom contained within the suite.

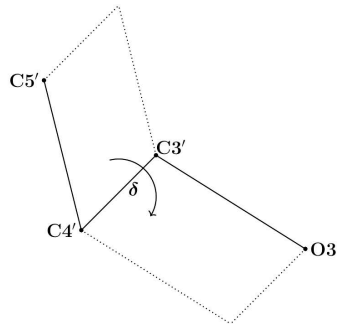
Viewed from the suite perspective, the RNA backbone has the attractive property of rotamericity [11]. That means bond geometry and steric contacts constrain the backbone to lie in generally distinct conformational clusters. **Fig 3B** names the backbone dihedral angles used by the RNA Society's multi-lab RNA Ontology Consortium to define consensus suite conformers [12]. The number and letter labels denote the parts of those conformer names, such as **1a**, **5z**, **6n**, etc. At the time of that work, naive clustering algorithms failed at satisfactorily identifying and separating suite conformers, at least partly because 2/3 of the suites in the single A-form helical conformation overshadowed the rarer conformations. Therefore the 2008 conformers were first separated automatically, and then further conformer class boundaries were assigned by hand in a 7D graphics system developed for that purpose [12]. Additionally we tested whether close clusters should be



**Fig 3. Atomic structure and nomenclature for the two residues and the central sugar-to-sugar suite of an RNA dinucleotide.** The bases are simplified to single green lines. Backbone atoms other than carbon are marked with balls, red for O, blue for N, yellow for P, and with gray bonds out to H. Panel A: Names of heavy atoms used in describing suites. Panel B: The 7 dihedral angles of the suite are labeled in order as  $\delta$ ,  $\epsilon$ ,  $\zeta$ ,  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ . Additionally we show the dihedral angle  $\nu$  measuring the sugar pucker (see **Fig 4**). Number and letter refer to suite conformer names as, e.g., in **Section 2.2** and Table D in **S1 Text**.

<https://doi.org/10.1371/journal.pcbi.1014242.g003>

angle	atom bonds
$\alpha$	O3' - P - O5' - C5'
$\beta$	P - O5' - C5' - C4'
$\gamma$	O5' - C5' - C4' - C3'
$\delta$	C5' - C4' - C3' - O3'
$\epsilon$	C4' - C3' - O3' - P
$\zeta$	C3' - O3' - P - O5'



**Fig 4.** Left: the dihedral angle (first column), gives the rotation from the plane spanned by the last three atoms with respect to the plane spanned by the first three atoms, as depicted in the exemplary right panel (figure and table adapted from [31]), see Fig 3 The first 7 measure backbone dihedral angles, distinguished from  $\nu$  which measures the sugar pucker, see Fig 2A.

<https://doi.org/10.1371/journal.pcbi.1014242.g004>

separated or not by whether they were distinct or interchangeable in their local structural motif or functional role. Suites with conformations that do not match any of the consensus suite conformers are considered outliers for validation purposes and given the !! designation (pronounced “bang-bang”).

The original suite conformers were identified based on the complete set of 7 backbone dihedrals. However, these dihedrals are only reliably modeled at resolutions where all their member atoms are clear in the experimental map. Just as the Pperp criterion predicts the correct sugar pucker from minimal model information, and even in the presence of modeling errors, RNAprecis seeks to predict, as accurately as possible, correct rotameric suite conformers from minimal model information.

## 1.2 Prior work and background in geometric statistics

Statistical analysis of geometric objects such as suite conformers requires firstly their mathematical representation. For instance, aggregating the positions  $\mathbf{x}_i \in \mathbb{R}^3$  of their  $k$  atoms ( $i = 1, \dots, k$ ) results in a configuration matrix

$$(\mathbf{x}_1, \dots, \mathbf{x}_k) \in \mathbb{R}^{3 \times k}$$

where each column is called a *landmark* vector.

Secondly, it is necessary to mathematically account for the fact that a translated and rotated suite conformer represents an object of the same shape. Thus the *shape* of a configuration matrix is the *configuration equivalence class* of all configuration matrices that can be reached from a single configuration by applying a single rotation and a single translation to all of the columns. Developing statistical methods for the analysis of such and related concepts of shapes has picked up momentum since the middle of the last century and a broad overview is given by [13]. The following two representations of shape are elaborated in [14]

As detailed above, a favorable representation of rotamer shape at high detail, avoiding configuration equivalence classes altogether, is given by a dihedral angle representation. Indeed, while individual bond lengths and angles between bonds at atomic centers are fairly fixed (with standard deviations on the order  $10^{-2}\text{\AA}$  for lengths and  $\sim 2$  degrees for angles [15]), structural diversity is caused by rotation around bonds. These rotations, given by dihedral angles, can be measured in arc length on the unit circle,

$$\mathbb{T} := [0, 2\pi] / \sim$$

**Table 1. Selection of tuning parameters required for the AGE algorithm from [19] for the four different pucker-pair HD data sets: relative branching distance  $q$ , minimal cluster size  $\kappa$  (set to  $\kappa = 3$  for all four data sets in order to also identify clusters with very few elements) and the maximum outlier distance  $d_{\max}$  (given as the percentage of elements that are in a cluster with less than  $\kappa$  elements).**

dataset	$q$	$\kappa$	$d_{\max}$
HDP33	0.09	3	1.1%
HDP32	0.05	3	0.34%
HDP23	0.07	3	1.4%
HDP22	0.05	3	0.7%

<https://doi.org/10.1371/journal.pcbi.1014242.t001>

where “/  $\sim$ ” indicates formation of the quotient space (circle) from the interval by identifying the endpoint 0 with the other endpoint  $2\pi$ . Then, the shape of one suite at high detail, determined by 7 dihedral angles, see Figs 3 and 4, is given by an element of the 7-dimensional torus

$$\mathbb{T}^7 := ([0, 2\pi]/\sim)^7. \quad (1)$$

**MINT-AGE.** In order to statistically cluster torus data in an unsupervised way, in previous work we had developed MINT-AGE (Mode huntINg after Torus pca on Adaptive iterative cutting linkaGe trEes). Its building block, AGE, can be applied to data on any metric space. AGE first builds a linkage tree (e.g., average linkage), as introduced by [16] and then performs a specifically designed iterative adaptive clustering method inspired by [17,18], identifying first the lowest density cluster, then iteratively higher density clusters. The method is detailed in [19] and our specific tuning parameters are listed in Table 1. The building block MINT is then applied to every (pre-)cluster found by AGE. Each cluster is subjected to *torus PCA* from [20]. Circumventing the fundamental difficulties of straightforward PCA approaches on a torus, here the 7-dimensional torus is mapped to a 7-dimensional stratified sphere (where the stratification mimics the topology of the torus) and an adaption of the *principal backward nested subsphere* algorithm (PNS) from [21] (a variant of PCA for spheres), using torus distances rather than spherical distances, is applied to the data, which is then projected to its main one-dimensional small circular component. These projected data are subjected to circular mode hunting (e.g., [22,23]) where each statistically significant mode stands for a new subcluster. Originally, MINT used nonparametric circular mode hunting ([20] building on the linear version by [24]), where cluster boundaries had been placed at antimodes of suitably kernel smoothed data utilizing a similar approach to the WiZer (Wrapped Gaussian kernel smoothed derivative’s significant Zero crossings, from [23] building on [25]). In order to identify also very small subclusters, even of size 2, we use a parametric approach, namely a likelihood ratio test for the likelihoods of fitted unimodal and bimodal Gaussians, recently developed in [26]. This test can be performed iteratively to find several subclusters on a main one-dimensional small circular component.

**Multicentred coordinates for constrained shape.** At low detail, where variation of pseudo-bond angles and lengths contribute considerably to shape variability, we resort to a landmark-based shape representation. Once again, however, we avoid the abstract modeling by configuration equivalence classes and develop a *frame-based* method. Here, specific landmarks (i.e., specific prominent atom positions) are singled out and mapped to a standardized frame, thus uniquely determining rotation and translation of the other landmarks. Frame based methods have been successfully utilized as *Bookstein coordinates* [27], *Goodall-Mardia coordinates* [28] and also for *affine* and *projective shape* [29]. In all of these approaches, one landmark is subtracted from all others, thus becoming the origin. This role will be taken by the landmark describing the central PO<sub>4</sub> group. For data at low detail the rotation will be uniquely determined by rotating the pseudo-bonds from  $P$  to each of the two C1’ atoms into a shared coordinate plane. Since data at low detail also contains the glycosidic bonds of nearly fixed lengths between each C1’ and its neighboring N1/N9, the location of the former will be subtracted from the latter, leading to *multicentred coordinates for constrained shape* as proposed in [14]. Fréchet means

[30] of MINT-AGE clusters will be computed in multicentred coordinates and varying cluster breadth will be accounted for by a Gaussian likelihood in Bayesian inference.

## 2 Data and methods

### 2.1 Data sets used

**Sugar pucker-pair sets.** Since the ribose puckers are bimodal, a good start on automatically identifying full-detail backbone conformers is to sort conformers into 4 categories by each suite's two sugar puckers at each side. We call these the *pucker-pair* sets:

P33, P32, P23, P22,

where C3'-endo and C2'-endo have been shortened to just 3 and 2. The P33 group is much larger than the others, because it includes the dominant A-form double-helical conformation (**1a**).

At high detail, the sugar pucker is determined directly from atoms in the ribose ring, using the  $\nu$  dihedrals (see [Figs 2 and 4](#)), where  $\nu_1$  refers to the first ribose in the suite and  $\nu_2$  to the second. At low detail, the ribose ring atoms are not assumed to be reliable, and sugar pucker is instead determined using the Pperp criterion (see [Fig 2B](#)). Note that the P atom from the next sequential suite (with the index  $i + 1$ ) is required for determining the second sugar's Pperp distance and pucker conformation.

**High-detail gold standard.** The first training dataset, denoted by HD, comprises complete atomic positions of RNA suites obtained by high-resolution (<2.0Å) X-ray crystallography and CryoEM. It is derived from a preliminary version of a carefully curated RNA residue dataset, RNA2023 [32]. The 3.150 release of BGSU RNA 3d Hub (<http://rna.bgsu.edu/rna3dhub/nrlist>) classified RNA chains of similar function and nucleic acid sequence into *sequence equivalence classes*. (Note that sequence equivalence classes are an entirely distinct concept from the configuration equivalence classes introduced and used only in [Section 1.2](#); the overlap in terminology represents one of the challenges in interdisciplinary work.) Then the best chain is selected from each sequence equivalence class based on resolution, RSR, RSCC, Rfree, clashes, and completeness [33]. Chains from two high-resolution CryoEM ribosome structures, namely 8b0x (E. coli) [34] and 8a3d (human) [35], were added to this set of chains.

Since even high-resolution structures may contain local disorder or local modeling errors [36], the selected RNA chains were then filtered at the suite level. Suites with strong indicators of modeling errors, such as Pperp/pucker mismatches, steric clashes, and covalent bond geometry outliers have been removed from the dataset. Likewise, suites with insufficient support from experimental data, as measured by local 2mFo-dFc map values (or inclusion in cryoEM map envelope), local real-space correlation coefficients, or alternate conformations were removed. Additionally, suites involving a non-standard RNA base (i.e., other than A, U, C, or G), and suites identified as "!!" suite outliers by `phenix.suitename` validation were removed.

This quality-filtered data set serves as the *gold standard*. Since puckers are also recorded, HD comprises the following four datasets:

HDP33, HDP32, HDP23, HDP22.

**Low-detail** datasets comprise only positions of the large atomic agglomerations, namely the positions of the central P of the phosphate groups and the ends of the glycosidic bonds, i.e., the positions of the C1' and N1 or N9 atoms, respectively. For training purposes, from the *high-resolution* dataset a low-detail dataset denoted by LD has been obtained by taking the only corresponding 5 atom positions, yielding, four datasets featuring corresponding sugar puckers:

LDP33, LDP32, LDP23, LDP22.

The PDB files of the high-resolution dataset are listed with the corresponding resolution in Table A in [S1 Text](#) and the HD and LD representations, with their pucker-pairs, are plotted in Figs A and B in [S1 Text](#).

**Low-detail test data and related conformer classes as best-guess answers.** For testing how well our trained system can predict correct conformations, test sets were created. These test sets address two related questions. First, are the model features used by RNAprecis reliably preserved in typical RNA modeling errors? Second, given reliable parameters, can RNAprecis predict best guess solutions for RNA modeling errors? The first test set, named LDT, contains only conformational outlier suites from structures solved at lower resolution than the training set. The second data set stems from the 23S large ribosomal subunit described by the PDB identifier 8b0x chain a (lower case), which is a member of the final RNA2023 dataset but has no correspondence in the training data and can therefore be considered an out-of-sample test set. These datasets comprise, for each member suite, measurements of the same five atomic positions as in LD (the central P and the two neighboring C1' and N1/N9), plus an additional sixth atom (P from the next residue following the suite). LDT is divided using the Pperp criterion into the following pucker-pair datasets:

LDTP33, LDTP32, LDTP23, LDTP22.

and correspondingly for the 8b0x dataset we have the pucker-pair datasets:

8b0xP33, 8b0xP32, 8b0xP23, 8b0xP22.

Suites for the test sets were selected to provide examples of modeling errors with probably known solutions from related structures. For each chain in the HD training set plus 8b0x chain a, related chains were identified based on the sequence equivalence classes from the 3.328 release of <http://rna.bgsu.edu/rna3dhub/nrlist> [33]. If a sequence equivalence class contained many chains, a representative chain was selected for the test set, based on that chain having an experimental resolution near the median resolution for that class and a count of!! outliers near the median count for the class. Such chains represent a “typical” degree of modeling challenge for their respective classes. A single equivalence class could contribute multiple chains to the test set if those chains had outliers at different sequence locations. Sequence equivalence classes of low sequence variety (such as poly-U) were excluded, as the relatedness of their member chains’ conformations is uncertain.

The chains in the test sets represent experimentally solved, PDB-deposited structures with resolutions from 1.39 to 3.1Å and release dates as recent as 2023.

A sequence alignment was performed on the selected chains with their reference chain from HD. From the aligned chains, suites that were modeled (i.e., identified) as!! outlier conformations were selected as members of the LDT test set if the reference chain had a recognized, non-!! conformer modeled at the same sequence location. The reference chain’s conformer was taken as the best guess correct answer. These reference chain suites all passed the stringent residue-level filtering for the RNA2023 dataset.

The PDB files, the chain names, and the residue numbers of the elements that comprise the low-detail test datasets are listed with their associated resolutions in [S1 Text](#). Table B in [S1 Text](#) lists the suites of the LDT dataset and the four pucker-pair sets are plotted in Fig C in [S1 Text](#), while Table C and Fig D in [S1 Text](#) provide the equivalent information for the 8b0x dataset.

Notably, we are not trying to analyze low detail suites where our presumptions are not met, i.e., where phosphates or glycosidic bonds are not placed well in density or do not have any density.

## 2.2 Conformer classes, MINT-AGE clusters and RNAprecis

**Conformer classes.** For each suite in the high-detail training data HD, the suite conformer class was determined with `phenix.suitename`. `phenix.suitename` uses the complete 7 dihedral angles of the suite ([Fig 3](#)) to assign the suite to

one of the manually defined rotameric conformers (e.g., **1a**, **1c**, **5z**, etc.) or to **!!** if the suite does not match any recognized conformation. The conformer class assignments from `phenix.suiteName` represent the current (best available to date, to best knowledge of the authors) “ground truth” suite identities for the HD datasets [12].

Each suite conformer class corresponds to a single pucker configuration pair—identifiable from the  $\nu$  dihedral angle at high detail—and training data suites were assigned to pucker-pair sets accordingly. The HD and LD datasets comprise 4148 suites that qualified for analysis (3529 are P33; of those, 3183 are **1a** and 346 non-**1a**; 216 are P23; 265 are P32; and 138 are P22.)

**Clustering** of the four high-detail datasets has been performed based on the purely statistical method of MINT-AGE developed by [31]. In order to accommodate very low conformer class sizes, we have used the new *mode hunting method for very small dataset sizes* from [14]. Tuning parameters have been learned such that the *clusters* found are in high agreement with the above suite conformers, see [Section 3](#).

**MINT-AGE training** was carried out on the HD datasets. The HD datasets are based on a preliminary version of the ‘nosuiteout’ RNA2023, as described above and available at <https://zenodo.org/doi/10.5281/zenodo.8103013> [32]. The files that have been used are listed in Table A in [S1 Text](#). The source PDB files were downloaded from the RCSB PDB database [37]. MINT-AGE used the full 7 dihedral angles to automatically cluster suites into clusters similar to the manually-defined suiteName conformer classes.

**RNAprecis training** was carried out on the LD datasets. The LD datasets comprise a minimal set of atom positions for the same files and suites as the HD datasets. The LD datasets are split into pucker-pair sets based on the Pperp ribose pucker criterion described above in [Section 1.1](#), which requires only atom positions from this minimal set. Because the training data is filtered to remove any ambiguous or incorrect ribose puckers, the pucker-pair sets are also reliably identified by  $\nu$  values.

**RNAprecis testing** was carried out on the LDT and 8b0x datasets. LDT datasets were based on the 3.328 version of the <http://rna.bgsu.edu/rna3dhub/nrlist> RNA list, dated 27.3.2024; files were downloaded from the RCSB PDB on 1.4.2024. The high-detail atom positions needed to determine  $\nu$  angles were not considered reliable for test set suites, so the suites were separated into pucker-pair sets based on the low-detail Pperp criteria ([Fig 8](#)). The LDT test set comprises 244 suites: 136 P33, 38 P32, 50 P23, and 20 P22. The 8b0x test set comprises 178 suites: 98 P33, 40 P32, 26 P23, and 14 P22. The PDB files and residues used are listed in Tables B and C in [S1 Text](#).

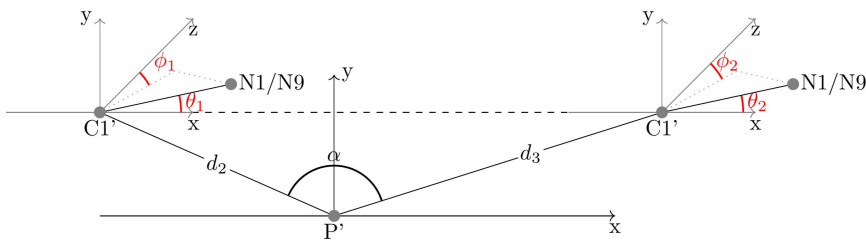
## 2.3 Non-Euclidean statistical learning methods for classification

**2.3.1 Multicentered coordinates for low-detail constrained size-and-shape.** Low-detail data is given by 5 landmarks

$$\mathbf{x}_1, \dots, \mathbf{x}_5 \in \mathbb{R}^3 \quad (2)$$

of atom center positions in three dimensions. Here the 5 landmarks denote the estimated atom centers of  $N1/N9_{i-1}$ ,  $C1'_{i-1}$ ,  $P'_i$ ,  $C1'_i$  and  $N1/N9_i$  in that order, see [Fig 3](#).

**Size-and-shape** information of such a configuration is given by filtering out Euclidean motions (translations and rotations), and to this end, as alluded to in the introduction, we single out a *frame* of prominent landmarks as illustrated in [Fig 5](#): The P atom center naturally serves as a reference point, i.e., its position is subtracted from all other atom centers. P in conjunction with the two  $C1'$  atom centers naturally defines a planar nondegenerate reference triangle, which we rotate into the  $(x,y)$ -plane such that the rays from  $P$  to each of the two  $C1'$  have the same projection onto the positive  $y$ -axis and such that the ray to the later-in-sequence  $C1'$  (see [Fig 3](#), usually but not always the longer of the two rays) has a positive projection onto the  $x$ -axis. Then, the opening angle  $\alpha$  at the P atom of the reference triangle lives on a proper half circle. The two landmarks giving the  $N1/N9$  atoms, however, are not centered with respect to the  $P'$  atom but with respect to



**Fig 5. MUCSSS coordinates for LD suite size-and-shape.** Part of the parameters describes the triangle spanned by the phosphorus as foot point and the two neighboring C1' atoms. After normalization for rotation and height, three descriptors remain: the lengths of the two pseudo-bonds and the opening angle. The rest of the parameters describe the two N1/9 atoms at the corresponding glycosidic bonds (which have fixed lengths) each expressed in spherical coordinates with respect to the x-y-plane defined by the triangle.

<https://doi.org/10.1371/journal.pcbi.1014242.g005>

their neighboring C1' atoms defining the glycosidic bond. Since, as noted in the introduction, these bonds lengths can be considered as a known constant, only the directions are variable, thus the coordinates of the two N1/9 atoms can be viewed to live on a unit sphere.

**Detailed derivation of coordinates.** As sketched above, first subtract the phosphorus' coordinates from the two C1' coordinates and the respective C1' coordinates from the corresponding N1/N9, to obtain the 4 new landmarks

$$\mathbf{y}_1 := \mathbf{x}_1 - \mathbf{x}_2, \mathbf{v} := \mathbf{x}_2 - \mathbf{x}_3, \mathbf{w} := \mathbf{x}_4 - \mathbf{x}_3, \mathbf{y}_2 := \mathbf{x}_5 - \mathbf{x}_4.$$

In order to rotate the triangle, which we may assume to be nondegenerate, spanned by the origin,  $\mathbf{v}$  and  $\mathbf{w}$  into the (x,y)-plane, as sketched above, define orthogonal unit vectors (where  $\mathbf{e}_1$  and  $\mathbf{e}_2$  define the new x- and y-axis, respectively)

$$\left. \begin{aligned} \mathbf{e}_1 &:= \frac{\mathbf{w}-\mathbf{v}}{\|\mathbf{w}-\mathbf{v}\|}, \\ \mathbf{e}_2 &:= \frac{\mathbf{v}-\mathbf{e}_1^T \mathbf{v}}{\|\mathbf{v}-\mathbf{e}_1^T \mathbf{v}\|}, \\ \mathbf{e}_3 &:= \mathbf{e}_1 \times \mathbf{e}_2, \end{aligned} \right\} \quad (3)$$

and the orthogonal matrix  $R := (\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3)$ . Indeed, then

$$\mathbf{e}_1^T \mathbf{v} = \frac{\mathbf{v}^T \mathbf{w} - \|\mathbf{v}\|^2}{\|\mathbf{w} - \mathbf{v}\|}, \quad \mathbf{e}_1^T \mathbf{w} = \frac{\|\mathbf{w}\|^2 - \mathbf{v}^T \mathbf{w}}{\|\mathbf{w} - \mathbf{v}\|}$$

and

$$0 < \mathbf{e}_2^T \mathbf{v} = \frac{\sqrt{\|\mathbf{v}\|^2 \|\mathbf{w}\|^2 - (\mathbf{v}^T \mathbf{w})^2}}{\|\mathbf{w} - \mathbf{v}\|} = \mathbf{e}_2^T \mathbf{w}.$$

Thus, the *frame-based, multicentered coordinates* for the above 5-point configuration (2) are given by

$$\left( \|\mathbf{v}\|, \|\mathbf{w}\|, \arccos \frac{\mathbf{v}^T \mathbf{w}}{\|\mathbf{v}\| \cdot \|\mathbf{w}\|}, R^T \mathbf{y}_1, R^T \mathbf{y}_2 \right) \in \mathbb{R}_+ \times \mathbb{R}_+ \times (0, \pi) \times \mathbb{R}^3 \times \mathbb{R}^3.$$

Taking into account that the two glycosidic bonds are constrained to have equal lengths that are the same for all configurations, we may assume  $\|\mathbf{y}_1\| = 1 = \|\mathbf{y}_2\|$ , to obtain the MUCCSS (MULTiCentred Constrained Size-and-Shape) coordinates given by

$$W := (d_2, d_3, \alpha, \mathbf{s}_1, \mathbf{s}_2) \in \mathbb{W} := \mathbb{R}_+^2 \times \mathbb{S}_+^1 \times (\mathbb{S}^2)^2, \quad (4)$$

where

- $\mathbb{S}_+^1 = (0, \pi)$  is the unit half circle,
- $d_2 := \|\mathbf{v}\|$  and  $d_3 := \|\mathbf{w}\|$  denote the distances between the atomic centers of P and the two C1', respectively,
- $\alpha := \arccos \frac{\mathbf{v}^T \mathbf{w}}{\|\mathbf{v}\| \|\mathbf{w}\|}$  gives the angle at P between the two pseudo-bonds from P to the two C1', respectively, and
- $\mathbf{s}_i := R^T \mathbf{y}_i$ ,  $i=1,2$  give the directions of the two glycosidic bonds, respectively, with polar angles  $\phi_1, \theta_1$  and  $\phi_2, \theta_2$ , respectively (see Fig 5),  $\theta_i \in [0, \pi]$ ,  $\phi_i \in [0, 2\pi]$ ,  $i=1,2$ .

Notably, MUCCSS coordinates for LD data have 7 degrees of freedom, this is also the number of degrees of freedom for HD data modeled on the 7-dimensional torus, see (1).

**Metric, means, tangent space coordinates and tangent space covariance.** Since the space of MUCCSS coordinates for LD data, given by (4), is a direct product, its canonical metric is that of a direct product:

$$d_{\mathbb{W}}(W^{(1)}, W^{(2)}) := \left( \sum_{j=1}^{n_j} \left( (d_2^{(1)} - d_2^{(2)})^2 + (d_3^{(1)} - d_3^{(2)})^2 + (\alpha^{(1)} - \alpha^{(2)})^2 + \arccos^2((\mathbf{s}_1^{(1)})^T \mathbf{s}_1^{(2)}) + \arccos^2((\mathbf{s}_2^{(1)})^T \mathbf{s}_2^{(2)}) \right) \right)^{\frac{1}{2}}, \quad (5)$$

for

$$W^{(j)} = (d_2^{(j)}, d_3^{(j)}, \alpha^{(j)}, \mathbf{s}_1^{(j)}, \mathbf{s}_2^{(j)}) \in \mathbb{W}, \quad j = 1, 2.$$

Further, for a sample

$$W^{(j)} = (d_2^{(j)}, d_3^{(j)}, \alpha^{(j)}, \mathbf{s}_1^{(j)}, \mathbf{s}_2^{(j)}) \in \mathbb{W}, \quad j = 1 \dots, n$$

of suite representatives in LD we have their *Fréchet function*

$$\mathbb{W} \rightarrow \mathbb{R}, \quad W \mapsto F(W) = \frac{1}{n} \sum_{j=1}^n d_{\mathbb{W}}(W, W^{(j)})^2$$

and every minimizer

$$W^* \in \arg \min_{W \in \mathbb{W}} F(W) \quad (6)$$

is a *Fréchet mean*.

Then, every  $W = (d_2, d_3, \alpha, \mathbf{s}_1, \mathbf{s}_2) \in \mathbb{W}$  has the following *tangent space coordinates* in the tangent space of  $\mathbb{W}$  at  $W^* := (d_2^*, d_3^*, \alpha^*, \mathbf{s}_1^*, \mathbf{s}_2^*)$ , given by

$$T_{W^*}(W) := (d_2 - d_2^*, d_3 - d_3^*, \alpha - \alpha^*, \tilde{\mathbf{s}}_1, \tilde{\mathbf{s}}_2)^T \in \mathbb{R}^7.$$

Here,  $\tilde{\mathbf{s}}_i$ ,  $i=1,2$ , is obtained as follows: With the unit vectors from (3) select

$$\mathbf{a} := \arg \min_{\mathbf{e}_k \in \{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}} |\mathbf{e}_k^T \mathbf{s}_i| \text{ and set } \mathbf{b} := \frac{\mathbf{a} - \mathbf{s}_i \mathbf{a}^T \mathbf{s}_i}{\|\mathbf{a} - \mathbf{s}_i \mathbf{a}^T \mathbf{s}_i\|}, \quad \mathbf{c} := \mathbf{b} \times \mathbf{s}_i, \quad \mathbf{Q} := (\mathbf{b} | \mathbf{c}),$$

to obtain

$$\tilde{\mathbf{s}}_i := \mathbf{Q}^T \mathbf{s}_i \in \mathbb{R}^2.$$

This yields the orthogonal projections of  $\mathbf{s}_i$  ( $i=1,2$ ) to the tangent space of  $\mathbb{S}^2$  at  $\mathbf{s}_i$  in a coordinate system only depending on  $\mathbf{s}_i$ .

Notably, for numerical robustness,  $\mathbf{b}$  and  $\mathbf{c}$  have been chosen from  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$  such that  $\mathbf{s}_i$  is mainly explained by them. Note also that the mapping  $\mathbf{s}_i \mapsto \tilde{\mathbf{s}}_i$  is only bijective if  $\mathbf{s}_i$  is in the hemisphere  $\{\mathbf{s} \in \mathbb{S}^2 : \mathbf{s}^T \mathbf{s}_i \geq 0\}$ . The data driven choice of  $\mathbf{b}$  and  $\mathbf{c}$  results in a robustification of the well known as *residual coordinates* [38].

Finally, the tangent space coordinates of a sample

$$W^{(j)} = (d_2^{(j)}, d_3^{(j)}, \alpha^{(j)}, \mathbf{s}_1^{(j)}, \mathbf{s}_2^{(j)}) \in \mathbb{W}, \quad j = 1 \dots, n$$

at a Fréchet mean  $W^* := (d_2^*, d_3^*, \alpha^*, \mathbf{s}_1^*, \mathbf{s}_2^*)$ , allow for the definition of a *tangent space covariance matrix* given by

$$S[W^{(1)}, \dots, W^{(n)}] := \frac{1}{n-1} \sum_{j=1}^n T_{W^*}(W^{(j)}) \left( T_{W^*}(W^{(j)}) \right)^T \in \mathbb{R}^{7 \times 7}. \tag{7}$$

*Remark 2.1 (Uniqueness).* The RNAprecis algorithm below requires Fréchet means of MINT-AGE clusters in  $\mathbb{W}$  representation. Due to [39], such sample means are unique on a complete connected Riemannian manifold if their distribution is absolutely continuous with respect to Riemannian measure. Indeed,  $\mathbb{W}$  is a connected Riemannian manifold. Although  $\mathbb{W}$  is noncomplete, as the boundary points corresponding to  $\alpha = 0, \pi$  are missing, it is geodesically convex, since the standard Riemannian product metric inducing the metric (5) is Euclidean on the problematic  $\mathbb{S}_+^1$ . Hence, these boundary points will never be coordinates of Fréchet means. In consequence, we can assume uniqueness of means and tangent space coordinates, and hence, also of tangent space covariance matrices.

**2.3.2 RNAprecis.** The following algorithm transforms LDT data to MINT-AGE clusters (introduced in Section 2.2) of HD gold standard data (see Section 2.1) and via its correspondence with suite conformation classes to suite conformations (see Fig 1).

**Algorithm 2.2 (RNAprecis). Input:**

1. Clusters obtained by MINT-AGE from HD gold standard data

$$C_1^{(p,q)}, \dots, C_{m(p,q)}^{(p,q)}, \quad m(p,q) \in \mathbb{N},$$

for every pucker configuration pair  $(p, q) \in \{2, 3\}^2$ ,

2. for every cluster  $C_i^{(p,q)}$ ,  $1 \leq i \leq m_{(p,q)}$  a Fréchet mean  $\bar{W}_i^{p,q} \in \mathbb{W}$  from (6) of its LD representatives along with the corresponding tangent space covariance matrix  $S_i^{p,q} \in \mathbb{R}^{7 \times 7}$  from (7),
3. a regularization parameter  $0 < \lambda < 1$  (with default value  $\lambda = 0.5$ ) and for all  $(p, q) \in \{0, 1\}^2$ ,  $1 \leq i \leq m_{(p,q)}$  the regularized covariance matrix

$$RS_i^{p,q} := \lambda S_i^{p,q} + (1 - \lambda) S_1^{3,3},$$

4. a test suite in LD representation

$$\bar{W} = (d_2, d_3, \alpha, \mathbf{s}_1, \mathbf{s}_2) \in \mathbb{W}$$

from the LDT test data.

**Steps:**

1. Determine the sugar pucker configuration pair  $(p, q) \in \{2, 3\}^2$  of  $\bar{W}$  via the Pperp criterion (see [Section 1.1](#)),
2. using cluster sizes  $|C_i^{(p,q)}|$  as (unnormalized) prior probabilities, determine (unnormalized, logarithmized) posterior probabilities

$$\log \tilde{P}_{\text{post}}(C_i^{(p,q)} | \bar{W}) := \log |C_i^{(p,q)}| - \frac{1}{2} \log \det (RS_i^{p,q})$$

$$- \frac{1}{2} T_{\bar{W}_i^{p,q}(\bar{W})} (RS_i^{p,q})^{-1} (T_{\bar{W}_i^{p,q}(\bar{W})})^T$$

and normalize them

$$P_{\text{post}}(C_i^{(p,q)} | \bar{W}) := \frac{\tilde{P}_{\text{post}}(C_i^{(p,q)} | \bar{W})}{\sum_{j=1}^{m_{(p,q)}} \tilde{P}_{\text{post}}(C_j^{(p,q)} | \bar{W})};$$

3. determine the most probable cluster index

$$i^* \in \arg \max_{1 \leq i \leq m_{(p,q)}} P_{\text{post}}(C_i^{(p,q)} | \bar{W}).$$

**Output:**

The conformation class(es) corresponding to the MINT-AGE cluster  $C_{i^*}^{(p,q)}$ .

*Remark 2.3.*

- 1) RNAprecis employs a Gaussian likelihood with estimated covariance for each cluster in order to ensure that variation of all coordinate representatives is taken equally into account: the *Mahalanobis distance* [40],

$$\tilde{d}_{\mathbb{W}}(W, W_i^{p,q}) := \sqrt{T_{W_i^{p,q}(W)} (RS_i^{p,q})^{-1} (T_{W_i^{p,q}(W)})^T}, \tag{8}$$

is the only term in the posterior probability which depends on the test suite coordinates. Then for fixed  $\mathbb{W}$ -distances between a test suite and a cluster mean, the more concentrated the cluster, the larger the distance  $\tilde{d}_{\mathbb{W}}$ .

- 2) Regularization for the covariance with the overwhelming largest cluster (**1a** has over 3,000, the second largest **1c** has only approx. 150 members) guarantees invertibility also for low-size clusters whose variation does not cover the entire 7-dimensional space.

3) The interpretation of posterior probabilities must take into account the covariance regularization. While these probabilities are valid posterior probabilities, they are connected to a model with regularized covariances, giving “regularized” credibility regions.

### 3 Results

#### 3.1 Clustering of HD and corresponding LD regions

In this section, the cluster results of the adapted MINT-AGE algorithm (see [Section 2.2](#)) for the individual data sets HDP33, HDP32, HDP23 and HDP22 (see [Section 2.1](#)) are reported and analyzed. The AGE algorithm requires 3 tuning parameters (see [19] and [31]) listed in [Table 1](#) for the four different data sets. The resulting clusters are listed in [Table D](#) in [S1 Text](#). Notably, as the confusion matrices in [Figs I, J, K, and L](#) in [S1 Text](#) teach,

8 out of 11 clusters in the HDP33 data set,

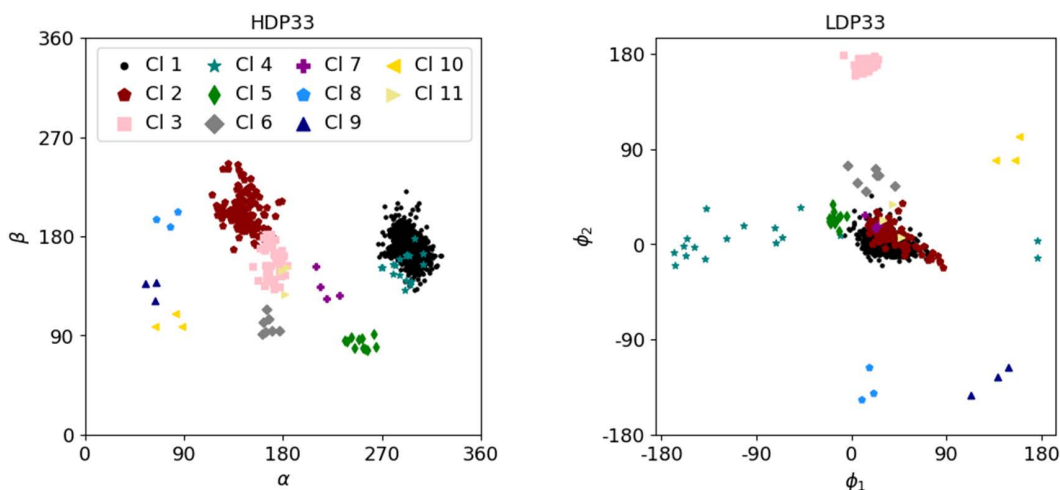
13 out of 15 clusters in the HDP32 data set,

10 out of 16 clusters in the HDP23 data set,

7 out of 10 clusters in the HDP22 data set

belong to a single suite conformer class. [Fig 6](#) shows the dihedral angle combination of  $\alpha$  and  $\beta$  for the data set HDP33 (by far the largest) in the left panel (see [Figs 3](#) and [4](#) for an explanation of the dihedral angles), in which the MINT-AGE clusters are best visually separated. The plots of all angle combinations for all data sets are shown in [Figs E, F, G, and H](#) in [S1 Text](#). There is also a detailed comparison between the cluster results and the suite conformer classes identified with `suitename`, using confusion matrices.

MINT-AGE successfully identifies clusters with small numbers of individual suites for many suite conformers. **7d, 5d, 5j, 4g, “0023a, 4d, 5p, 5q, 0b, 2u, and 2o** are all assigned unique clusters based on 4 or fewer examples. Nevertheless, some conformers remain unclustered: **7a, 3d, 5n, 3g, 5r, 0k, and 4s** could not be assigned to unique MINT-AGE clusters.



**Fig 6. Left panel: Illustration of HDP33 clusters in the  $\alpha - \beta$  plane of the high-detail coordinates.** Right panel: The low-detail coordinates corresponding to the same suites, representing the distribution of the clusters in the  $\phi_1 - \phi_2$ -plane. One can see that many of the clusters already well visible in the dihedral angle plot (left panel) are mapped to distinct clusters also in the low-detail plot (right panel) of longitudes of glycosidic bond directions.

<https://doi.org/10.1371/journal.pcbi.1014242.g006>

Conformers **7a**, **3d**, **5n**, and **3g** each have enough individual suites for MINT-AGE clustering, but their data clouds are obviously too spread out to assign to a single cluster, while **0k**, **5r**, and **4s** comprise only one individual suite each.

### 3.2 MINT-AGE clusters comprising multiple conformer classes

In broad overview (see Figs I, J, K, and L in [S1 Text](#)), HDP33 Clusters 1 and 4 combine two or more different conformer classes into a single cluster, as do HDP32 Clusters 1 and 5. It was expected that the corresponding suites would be difficult to separate automatically, as their conformer classes boundaries were based on structural context in the previous work. We now go into detail.

HDP33 Cluster 1 combines **1a**, **8a**, **1L**, and **1m**. **1a** is the dominant conformer class (A-form helix) across all RNA and exerts a strong influence on clustering. The **8a**, **1L**, and **1m** conformer classes are shoulders on the main **1a** cluster, each with a distinct behavior, and are not cleanly separated from **1a**. Similarly to the **1b/1l** case, these satellite conformers are distinguished from **1a** more by context than by conformation.

HDP33 Cluster 4 combines suite conformer Classes **9a**, **3a**, and possibly **7a**. These conformer classes differ only in the  $\zeta$  dihedral. As in the other cases, the distribution is quite continuous, and the distinction is found in the suite's context. **9a** conformers start or end loops, while **3a** conformers enable stack switches. **7a** conformers, of which a single instance is included in Cluster 4, are also stack switches.

HDP32 Cluster 1 combines suite conformer Classes **1b** and **1l**, as does Cluster 5. The **1b/1l** distinction is illustrative of the challenge that MINT-AGE faces in reproducing the manual conformer classes. In the `suitename` parameter space, **1b** and **1l** are closely related and share a continuous distribution with no clear conformational breakpoint between the behaviors. Instead the primary distinction comes from elsewhere in the structure: **1l** has a ligand or another base intercalated between the bases of the suite, while the bases of **1b** are often stacked directly on each other. In defining the original suite conformer classes [12], the cutoff between **1b** and **1l** was manually set at  $\beta=206^\circ$  to provide the cleanest possible separation between these behaviors.

Some clusters combine a main conformer with a small conformer that is otherwise unclustered. HDP33 Clusters 1 and 4 each included a single instance of **7a**. HDP23 Cluster 2 is majority **0a**, but also includes four **#a**. HDP23 Cluster 5 is majority **8d** and also contains two **4d**. HDP22 Cluster 6 does not contain a majority conformer, but combines three **2u** with two **2o**. These cases reflect the limitations of our relatively small training dataset; there are not enough examples of **#a**, **4d**, **2u**, and **2o** to create independent clusters for these conformers.

Several clusters include one or two instances of a conformer that is otherwise clearly identified in its own cluster. HDP23 Cluster 2 includes two instances of **4a**, which is otherwise cleanly assigned to Cluster 7. HDP23 Cluster 3 includes two stray **0i** and one stray **6j**. HDP23 Cluster 11 includes one stray **9d**. These cases are a reminder that RNA conformer space is complex and biological conformer boundary definitions are often somewhat ambiguous.

Finally, some conformers and clusters are particularly challenging to untangle. HDP23 Cluster 4 is majority **6g**, but contains two instances of **4g**, while **4g** is further split across Clusters 15 and 16. Cluster 16 also contains a single **4a**, which is otherwise assigned to Cluster 7. HDP22 Cluster 7 combines 4 instances of **4b** with two of **0b**, both of which conformers themselves are heavily split across multiple clusters. Expansion of the training dataset should help resolve these ambiguous conformers. However, the splitting of suite conformers into multiple clusters is not necessarily due to the limitations of the training data and may, as discussed below, present an opportunity to revisit the existing suite conformer definitions.

### 3.3 Conformer classes split across multiple MINT-AGE clusters

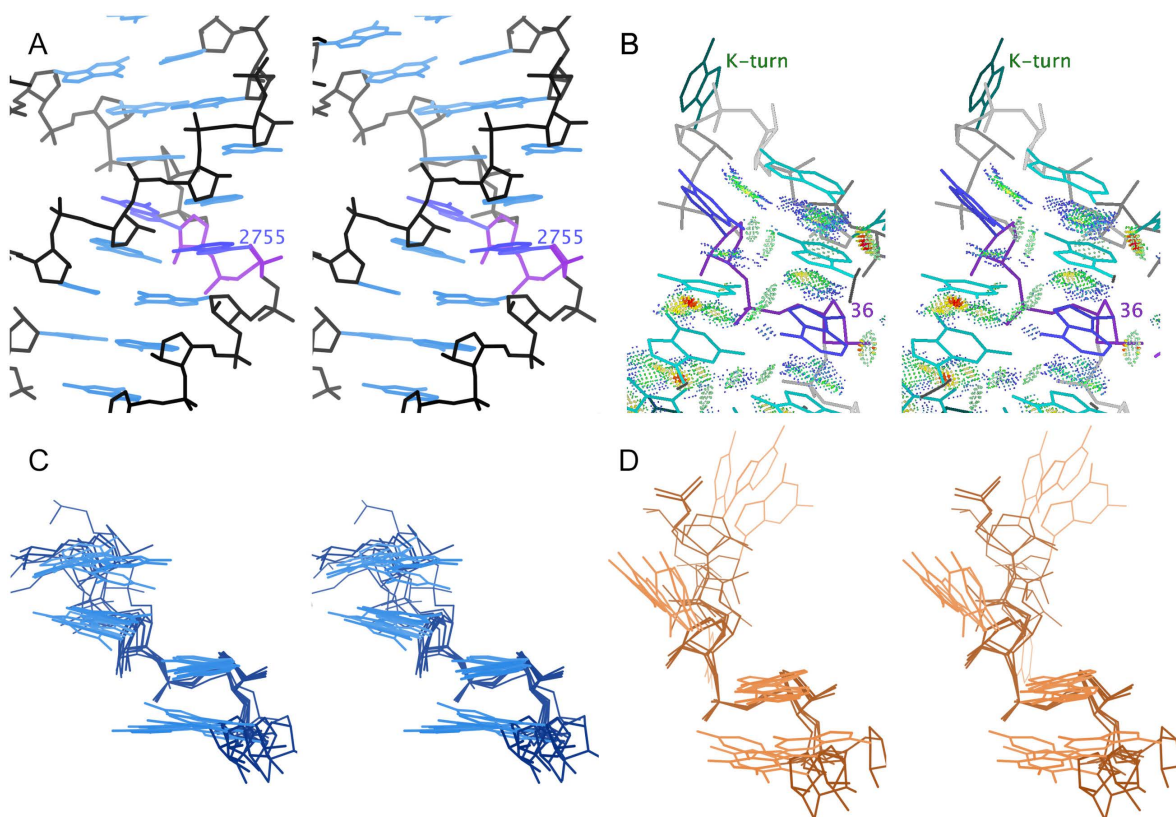
Some suite conformer classes were split into multiple clusters by MINT-AGE, most notably **0a** in HDP23, but also **6d** in HDP23, **7p** in HDP32, **1f** in HDP33, and **0b** in HDP22 (see Figs I, J, K, and L in [S1 Text](#)). For reproducing the suite conformer classes, it is easy to simply re-combine the clusters. However, the separation suggests that MINT-AGE has discovered real conformational divisions that were not visible during the previous work.

MINT-AGE strongly separates the **0a** conformer into two subclusters, with 16 suites in Cluster 2 and 12 suites in Cluster 6 (see Fig 7, and Table D in S1 Text). These subclusters separate the two major structural motifs found within **0a**, the stack switch and the kink-turn.

Cluster 2 contains stack switches. In standard **1a** conformation, the subsequent bases of the suite are near parallel and stack on top of each other. In these stack switches, the bases are still roughly parallel but change from base-stacking with each other to stacking with a different strand, usually in the same helix (Fig 7A and 7C).

Cluster 6 contains kink-turns. The kink-turn motif [42] spans multiple suites – the most distinctive feature is a very tight turn with an outturned base that often interacts with proteins. The Cluster 6 suite comprises the two residues following the residue with the outturned base. The first base of the Cluster 6 suite stacks with the backbone of the residue on the other side of the outturned base, stabilizing the sharp turn. The bases of a Cluster 6 suite are angled relative to each other due to this interaction with the preceding backbone (Fig 7B and 7D).

**0a** was previously treated as a single conformer class because these suites have a continuous distribution along the  $\zeta$  dihedral angle, best seen in the  $\zeta/\alpha$  panel of Fig F in S1 Text. However, Clusters 2 and 6 are more clearly distinguishable by  $\epsilon$ , which also separates the related conformer Class **4a**/Cluster 7, best seen in the  $\epsilon/\zeta$  panel of Fig F in S1 Text.

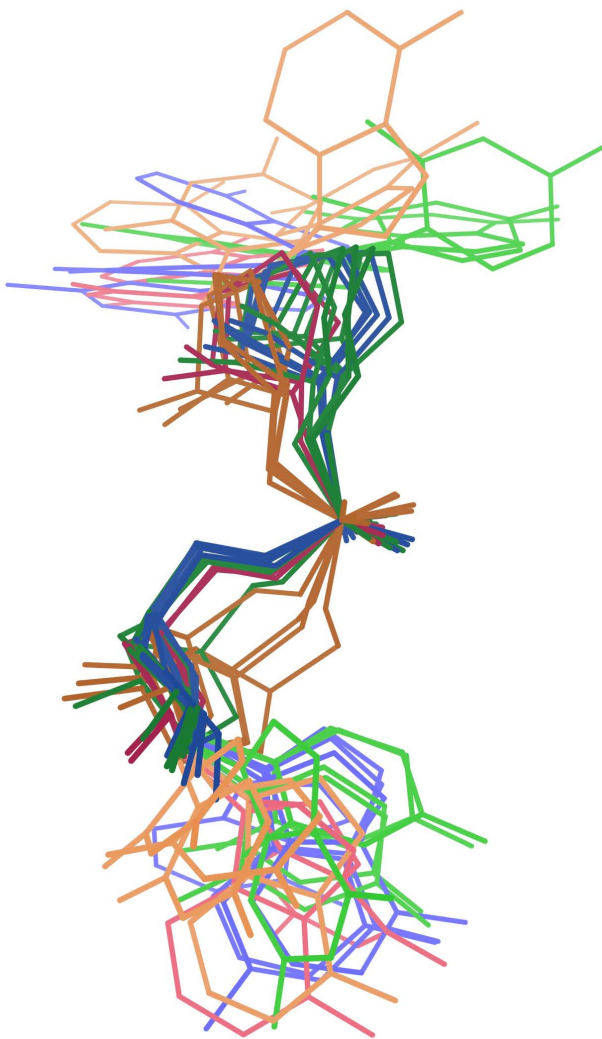


**Fig 7. Stereo images of **0a** conformers, as separated by MINT-AGE.** **A.** Stack switch from 8a3d chain A ([35]). Suite 2755 is highlighted in purple, and its bases stack into two different helix strands. **B.** Kink-turn from 5fjc chain A ([41]). Suite 36 is highlighted in purple. Below the outturned base, the upper base of suite 36 shows blue and green dots for favorable van der Waals packing with the backbone of the strand on the right. **C.** Superposition of 6 representative MINT-AGE Cluster 2 suites. These are stack switches and have bases close to parallel. **D.** Superposition of 4 representative MINT-AGE Cluster 6 suites. These are kink-turns and have bases at a high angle to each other.

<https://doi.org/10.1371/journal.pcbi.1014242.g007>

The split of suite conformer Class **0a** provides the strongest evidence of MINT-AGE identifying a systematic functional/conformational distinction that was not previously recognizable. The **7p** split in HDP32 is also suggestive. MINT-AGE finds 4 clusters, combining 7 instances into Cluster 9, 5 into Cluster 10, 4 into Cluster 13, and 3 into Cluster 14. This set of clusters is complex, but Cluster 9 is clearly conformationally distinct from the rest (Fig 8). That this distinction is conformational rather than “just” contextual demonstrates the value in approaching suite conformer identification from both perspectives.

The **0a** and **7p** suite conformer classes should be split into new conformer classes, based on the MINT-AGE clustering. The success of an unsupervised method in recognizing new conformer classes not previously identified by manual work is an encouraging result.



**Fig 8. Superposition of 7p suites, colored by MINT-AGE cluster.** Cluster 9 is orange, 10 is green, 13 is blue, 14 is red. Cluster 9 is clearly conformationally distinct from the other groupings within 7p.

<https://doi.org/10.1371/journal.pcbi.1014242.g008>

### 3.4 Statistical classification of conformations for low-detail data

The detailed results of the RNAprecis test are presented in [S1 Text](#), individually for the four pucker-pair data sets LDTP33 (Section E.1), LDTP32 (Section E.2), LDTP23 (Section E.3) and LDTP22 (Section E.4) formed by the classification in Step 1 of Algorithm 2.2. Analogously, 8b0xP33 (Section E.5), 8b0xP32 (Section E.6), 8b0xP23 (Section E.7) and 8b0xP22 (Section E.8) are treated separately. For each of the test elements, the probable answer conformer class determined by related experimental structures is provided for comparison (see paragraph “Low-detail test data and related conformer classes as probable answers” in [Section 2.1](#)). There are three possibilities, whose colors correspond to tables and figures below:

**Prediction matches:** The cluster determined by RNAprecis (Algorithm 2.2) contains a suite conformer class that matches the sequence equivalence-related probable answer for this element (see Table D in [S1 Text](#)). Note that clusters containing multiple conformer classes, such as P32 Cluster 1 which contains both **1b** and **1i**, are considered a match for any of their member conformers.

**Prediction mismatches:** The cluster determined by RNAprecis (Algorithm 2.2) contains no members matching the probable answer for the test element (see Table D in [S1 Text](#)). However, the probable answer conformer is a member of the same pucker-pair set as the test element.

**Pucker mismatches:** The pucker configuration pair of the cluster determined by RNAprecis from the test element (in Step 1 of Algorithm 2.2) does not match the pucker-pair set of the probable answer.

**3.4.1 Results for LDT test set.** For all four pucker-pair sets except P23, at least 2/3 of the elements are sorted into the prediction match group (blue), see [Table 2](#) and top left in Figs Q, R, T, and V in [S1 Text](#). The lower success rate for the P23 pucker-pair set is not clearly associated to a specific type of misclassification. However, in the P33 pucker-pair set, suites with probable answer **1a** are in large numbers classified as **1c**, a systematic effect discussed below. The probable sequence-related conformer class answers and the clusters determined with RNAprecis for all individual elements (sorted according to the four pucker-pair sets and according to the blue, red and orange categories) are reported in Table E (blue category of LDTP33), Table F (red and orange category of LDTP33), Table G (blue, red and orange category of LDTP32), Table H (blue, red and orange category of LDTP23), and Table I (blue, red and orange category of LDTP22), all in [S1 Text](#). A comparison between the probable answers and the clusters determined by RNAprecis is shown in a classification matrix for the four data sets in [Fig 9](#) below and in Figs S, U, and W in [S1 Text](#). The prediction results are separated into the four pucker configuration pairs so that the **1a** conformer does not dominate the results. The RNAprecis prediction results are good uniformly over all pucker-pair sets and all clusters, even small clusters.

The conformers **3d**, **5n**, **3g**, **0k**, **5r**, and **4s** were not clustered by MINT-AGE, probably due to having too few examples in the training dataset (between 1 and 4). Therefore, these conformers also could not be predicted by RNAprecis. Any test set suite where the probable answer conformation is one of these unclustered conformers will necessarily result in a prediction mismatch or pucker mismatch. This inability to predict matching conformers affects 3 suites in the test set, one each of **3d**, **3g**, and **5n**.

**3.4.2 Results for 8b0x test set.** In the “out of sample” test set from PDB 8b0x, RNAprecis again achieves high classification rates across the board, see [Table 3](#). However, in this case the performance on the P33 pucker-pair set lags

**Table 2. List of the four pucker-pair test datasets in the first column formed by Step 1 in RNAprecis (Algorithm 2.2) with the respective total counts in the second column and counts and percentages in the third through fifth columns.**

dataset	count	prediction matches	prediction mismatches	pucker mismatches
LDTP33	136	91 (67%)	37 (27%)	8 (6%)
LDTP32	38	30 (79%)	3 (8%)	5 (13%)
LDTP23	50	28 (56%)	15 (30%)	7 (14%)
LDTP22	20	16 (80%)	2 (10%)	2 (10%)

<https://doi.org/10.1371/journal.pcbi.1014242.t002>

		LDTP33															
RNAprecis	CI 6(1f)	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
	CI 9(5d)	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0
	CI 7(1e, 1f)	1	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0
	CI 4(9a, 3a, 7a)	1	0	0	3	0	0	0	0	0	1	1	1	0	0	0	0
	CI 3(1g)	1	0	8	0	0	0	0	0	0	0	0	0	0	0	0	1
	CI 1(1a, 1L, &a, 1m, 7a)	32	2	0	0	0	2	0	1	0	0	0	0	1	0	1	0
	CI 2(1c)	24	43	0	0	1	1	2	1	0	0	0	0	0	1	0	0
		1a	1c	1g	2a	1f	1b	1e	1L	5d	3a	3g	5n	0a	1l	1m	3d
		possible answers															

**Fig 9. Classification matrix listing counts of probable sequence-related conformation classes (horizontal axis) versus MINT-AGE clusters (vertical, with conformer classes in parentheses from MINT-AGE training, see Fig 1 in S1 Text) for the low-detail sugar test pucker-pair data set LDTP33 with colors: blue for match, orange for mismatch and red for pucker-pair mismatch.** Columns listing multiple suite conformers are from sites solved in different conformations across sequence-related models, and any of the sequence-related conformations were accepted as a matching predictions.

<https://doi.org/10.1371/journal.pcbi.1014242.g009>

**Table 3. List of the four pucker-pair test datasets for the 8b0x test data in the first column formed by Step 1 in RNAprecis (Algorithm 2.2) with the respective total counts in the second column and counts and percentages in the third through fifth columns.**

dataset	count	prediction matches	prediction mismatches	pucker mismatches
8b0xP33	98	53 (54%)	32 (33%)	13 (13%)
8b0xP32	40	33 (82%)	4 (10%)	3 (8%)
8b0xP23	26	20 (77%)	2 (8%)	4 (15%)
8b0xP22	14	9 (64%)	5 (36%)	0 (0%)

<https://doi.org/10.1371/journal.pcbi.1014242.t003>

behind due to the **1a** to **1c** misclassification also found in the LDT test set. Pucker mismatches are rare also in this test set. Notably, the 8b0xP23 configuration exhibits a lower rate of mismatches than the LDTP23, which further underscores that no systematic pattern of misclassification is apparent. Detailed results are presented in Sections E.5, E.6, E.7, and E.8 in S1 Text.

Conformers not clustered by MINT-AGE were more prevalent in the “out of sample” test set. 9 prediction mismatches, including 3 **3d** suites, 1 **3g**, 2 **4s**, 1 **5n**, and 2 **5r** are attributable to unclustered conformers. This accounts for the worse apparent performance in the 8b0xP22 test, which contains the 2 **4s** suites, while LDTP22 contained no unclusterable suites.

While 8b0x a has provided a useful test for demonstrating that RNAprecis does not overfit, we expect that including this chain in future MINT-AGE training will allow clustering and prediction of additional suite conformers. In particular, we expect to be able to identify **4s**, which occurs as a well-defined part of the “S-motif” [43].

## 4 Discussion

### 4.1 MINT-AGE clustering

This clustering process is quite successful in both training and test stages. In the training stage with filtered, high-detail data, 38 of the 51 suite conformers with  $\geq 3$  individual suites (75%) were clustered with one-to-one agreement by

MINT-AGE and only 4 were not clustered at all. The rest either combined 2 or more suite conformers into a single cluster or split one conformer into more than one cluster. Interesting and perhaps correct examples of those cases are described in [Section 3](#) (Results).

For the testing stage, RNAprecis uses the trained cluster relationships in reverse to predict the high-detail suite conformer for each!! suite that has a probable valid conformer (or several) in related structures. 68% were clear wins of prediction matches to a sequence-related conformation, and 23% were prediction mismatches. The other 9% were mismatches of the Pperp-determined pucker-pair set (used by RNAprecis) to the pucker-pair set of the probable answers, indicating cases where the modeling errors were too severe for this method. The details of the two mismatch categories have taught us several feasible ways to make the RNAprecis process significantly better, described below.

The most obvious way to improve MINT-AGE performance is use more training data. The RNA2023 dataset, from which the HD and LD training sets were adapted, is a general purpose dataset. One of its purposes is to reflect the relative frequency of different features in RNA structures. This means that several rare suite conformers (e.g., **3d**, **5n**, **4s**) are too rare (see above) in the dataset to form clusters with MINT-AGE. A new dataset, built to include more conformational diversity and more rare suite conformers (for instance by allowing two entries per sequence equivalence class) would improve clustering and recognition of those conformers.

MINT-AGE clustering can also be improved through development of additional parameters. For instance, the N1/N9-to-N1/N9 distance across the suite may differentiate between stacking (shorter distance) and intercalation (longer distance) in **1b** and **1j**. Such parameters would be targeted to specific conformer families and invoked only when needed.

## 4.2 Pucker mismatches

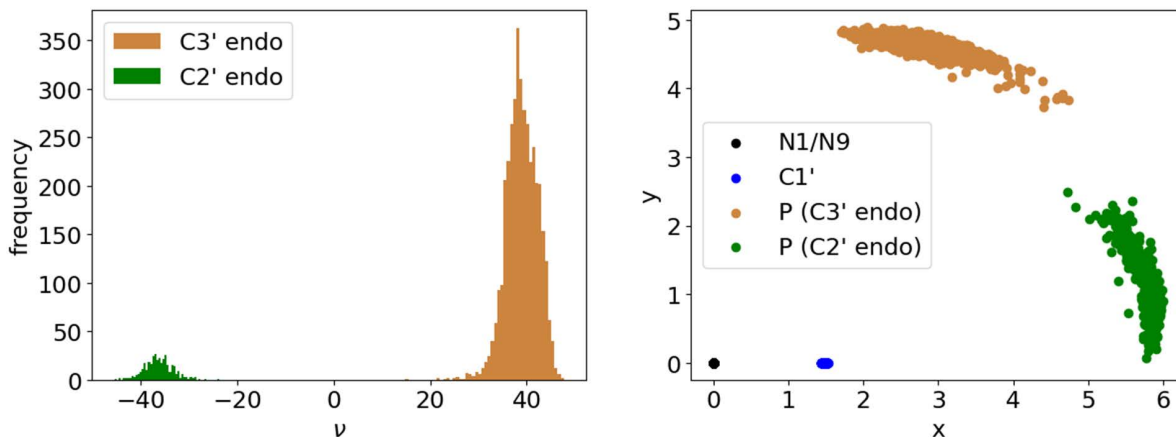
The most concerning category of incorrect predictions are the pucker mismatches, although they account for only about 9% of the total predictions. In these cases, the ribose pucker combination predicted from a test suite's Pperp distance test does not match the pucker configuration pair of the sequence equivalence class-related answer conformation. For example, 2qbz, chain X, suite 47–48 has a Pperp-predicted sugar pucker configuration pair corresponding to P23. The probable answer for this site (**6p**) has a different pucker configuration pair, corresponding to P22.

Quality-filtered residues from HD and LD are cleanly separated into pucker-pair sets by Pperp distance ([Fig 10](#)). Because the separation is so clear, MINT-AGE and RNAprecis treat each pucker-pair set separately, and RNAprecis prediction across pucker-pair sets is not possible. However, residues from the LDT test set can have ambiguous Pperp values near the 2.9Å cutoff ([Fig 11](#)). The otherwise useful Pperp criterion loses predictive power near this boundary.

A poor experimental map—or a poor fit to map—is the clearest cause of pucker mismatches. If the phosphate is not modeled near the center of a high map peak, or if the base and ribose are not each in a blob with plausible local relationships, then the basic assumptions of RNAprecis are not satisfied. [Fig 12](#) shows two such cases. Even without a pucker mismatch, RNAprecis is unlikely to find a correct answer in these cases, since the modeled suite does not necessarily bear any relationship to the correct conformation. Pucker mismatches thus represent a more severe category of modeling error than RNAprecis is intended to address, where the conformation of the model is very far from the intended structure.

The test set likely under-represents the actual prevalence of regions with experimental maps that cannot confidently support a model. All suites in the test set were solved as a recognized conformer in a related structure, so truly disordered regions were selected against. While Pperp distances in cases of pucker mismatch are considerably closer on average to the 2.9 boundary, than for correct matches (see [Fig 11](#)), the variation is very high such that this property is not suited to predict pucker mismatches. Other methods that directly assess map or fit-to-map quality will be necessary to identify suites that are likely to contain pucker mismatches. In regions with very poor experimental support, it is usually better simply to build no model at all, rather than try to predict conformations from compromised data.

An illustrative case for the challenges in RNA backbone prediction is 2qbz, chain X, suite 138–139. This suite is predicted with a pucker mismatch, however, the key features are all reasonably well-placed in the map. The problem is in the



**Fig 10.** For the training data of both high quality and high resolution, the two different sugar pucker can be very well discriminated by the angle  $\nu$  (left). Right: Using the coordinates described in Fig 2B, namely mapping N1/N9 to the origin and C1' to the first positive axis ( $x$ ), the perpendicular distance (Pperp, positive vertical  $y$  axis) similarly perfectly separates the two pucker configurations (this 2D display reflects true atom position in Euclidean space).

<https://doi.org/10.1371/journal.pcbi.1014242.g010>

preceding suite, 137–138. This suite has a severe error across the  $\alpha - \beta - \gamma$  dihedrals, which results in an incorrect sugar pucker in ribose 138. This pucker error interferes with the RNAprecis prediction for 138–139. This sort of compound error requires an iterative validation-remodeling-refinement process to correct step by step.

#### 4.3 1a vs 1c prediction mismatches

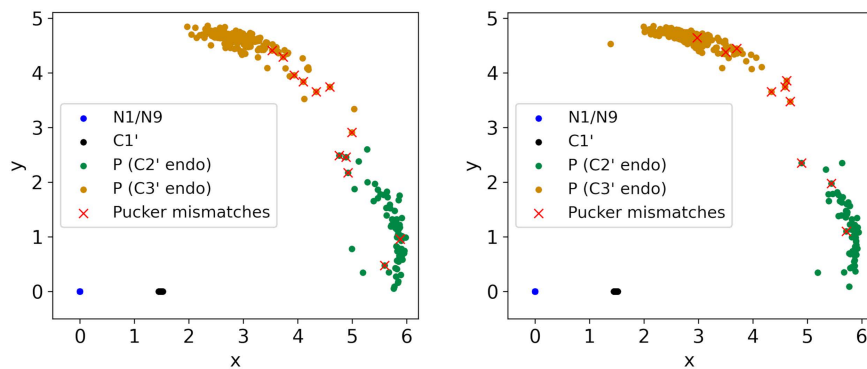
The classification matrix (Fig 9) for the LDTP33 portion of the test set shows that 1a vs 1c was one of the most difficult conformation pairs for RNAprecis to predict correctly. The difficulty is asymmetrical, with few 1c suites being misidentified as Cluster 1, but many 1a suites misidentified as Cluster 2.

A 1a/1c ambiguity is not unexpected. A coordinated motion of  $\alpha$  and  $\gamma$  allows  $\approx 180^\circ$  changes in those dihedral angles to occur with minimal effects on the parameters that RNAprecis uses (Fig 13). This motion resembles that of a physical “crankshaft”, having opposite and equal changes in  $\alpha$  and  $\gamma$  that drive changes in C5'/O5' atom positions, while  $\beta$  is held unchanged. In the Fig 13 superpositions, the glycosidic bond vectors of the lower base are systematically different in 1a versus 1c, but the difference is subtle relative to the huge size and population of the 1a conformer. RNAprecis therefore struggles to distinguish between these conformers.

We examined whether RNAprecis might be identifying the correct conformations in cases where the reference structure was mistaken. The reference structures of 1a/1c prediction mismatches were manually reviewed based on visual agreement with their experimental maps. The reference structure suites showed fits to map ranging from plausible to definitely correct. This manual review of the 1a/1c prediction mismatches also served to confirm that the sequence-related probable answers for the test set are reliable enough to use as targets. Where RNAprecis disagrees with the probable answers for this known difficult case, the sequence equivalence class-related conformation is generally correct.

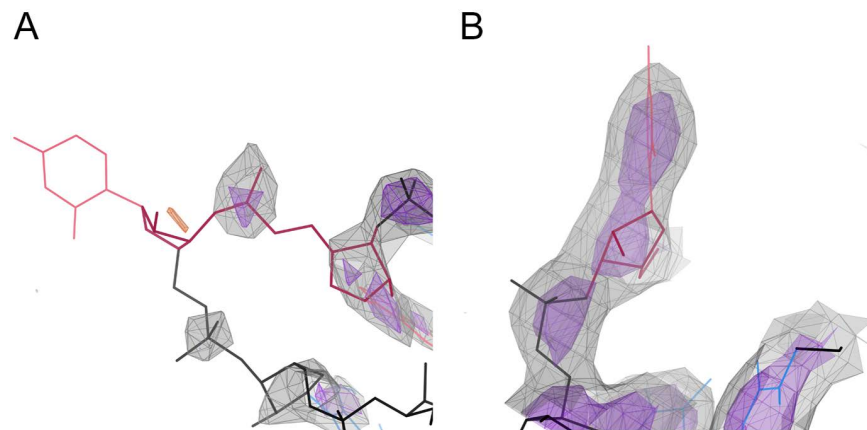
The 1a/1c ambiguity may be partially resolvable by increased reliance on the  $d_3$  parameter in this case.  $d_3$  is the distance across the relevant  $\alpha$ ,  $\beta$ , and  $\gamma$  dihedrals (see Fig 5). While there is some overlap between the Cluster 1/1a and Cluster 2/1c distributions in all parameters,  $d_3$  provides the cleanest single-parameter distinction between the two, as seen in the  $d_2/d_3$  panel of Fig M in S1 Text.

The asymmetry of 1a/1c misidentification affords a pragmatic immediate solution. Cluster 2 should be considered a possible prediction of either 1c or 1a, much as Cluster 1 is considered a match to 1a or any of its satellite conformers (but



**Fig 11. Depicting all low-detail test suites in the coordinates described in Fig 2B, namely mapping N1/N9 to the origin and C1' to the positive x-axis.** The left panel shows the respective coordinates for their first sugars (5' end of the suite), the right panel shows the respective coordinates of their second sugars (3' end of the suite). Considering  $y$  (the perpendicular distance to the positive  $x$ -axis), as in Fig 10, we predict the unknown true sugar pucker by vertically separating at 2.9, as learned from the high-resolution data, depicted in Fig 10. Matches of pucker predictions with probable answers are green and orange, respectively, red crosses indicate coordinates of pucker mismatches (see Section 3.4). There are only 12 pucker mismatches in the left panel and 10 in the right panel.

<https://doi.org/10.1371/journal.pcbi.1014242.g011>



**Fig 12. Two examples of problems causing RNAPrecis pucker mismatches.** A: Suite 47-48 of 2qbz, chain X, from LDT. Base 47 of this!! outlier suite has no density to define its position, and the corresponding ribose has a small difference density peak (orange,  $-3.5\sigma$ ) which may indicate a misplacement. Map contours at  $1\sigma$  (gray) and  $1.5\sigma$  (purple). B: Residue 561 of 7st2, chain 3. This base is part of!! suites 560-561 and 561-562, both from LDT. This base is clearly modeled in the wrong orientation relative to its map, resulting in an incorrect glycosidic bond vector. CryoEM map contours at  $6\sigma$  (gray) and  $10\sigma$  (purple).

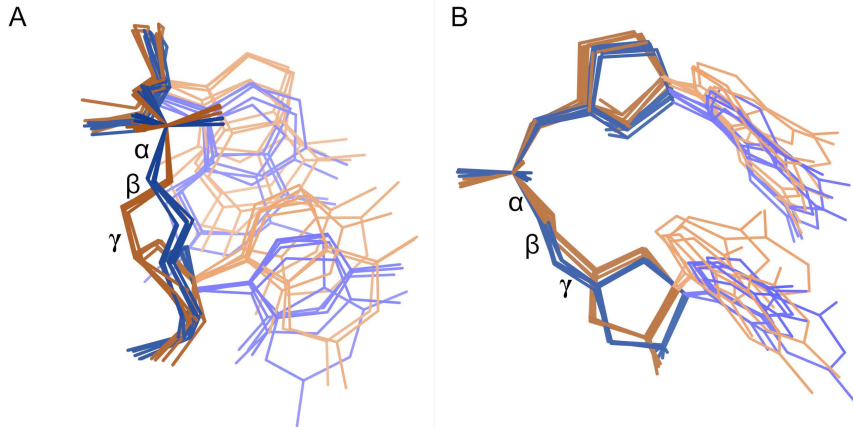
<https://doi.org/10.1371/journal.pcbi.1014242.g012>

not 1c) and Cluster 7 is considered a match to either 1e or 1f. In any of these cases, each member conformer of a cluster must be tested when making corrections to a structure.

#### 4.4 Other prediction mismatches

We manually examined all non-1a/1c prediction mismatches, comparing the test set suites to their experimental maps and to their corresponding probable answer conformations using the KiNG viewer. We did not find a unifying explanation for prediction mismatches, but some features of interest emerged.

Several prediction mismatch test suites from ribosome structures—including 7f5s chain L5 suite 4554, 8eiu chain A suite 195, and 7st2 chain 3 suites 149, 306, 519, 642, 1502, and 1519—showed conformations similar to their references



**Fig 13. Superposition of representative 1a (orange) and 1c (blue) suites.** **A.** View showing full backbone detail. **1a** and **1c** are easily distinguished by differences in the  $\alpha$  and  $\gamma$  dihedrals. Therefore, MINT-AGE can easily distinguish these conformers. **B.** Another view of the same superposition, showing the alignment of the riboses and bases. Despite the large difference in  $\alpha$  and  $\gamma$ , the related changes in glycosidic bond directions are relatively small. Therefore, RNAprecis cannot readily see the difference between **1a** and **1c**.

<https://doi.org/10.1371/journal.pcbi.1014242.g013>

answers except for a rotation of the  $\alpha/\gamma$  crankshaft. Given RNAprecis's difficulty resolving the **1a/1c** conformers, which differ by this same crankshaft motion, this collection of suites is suggestive of a systematic challenge. However, outside the ribosome examples, only one prediction mismatch suite, namely 3ger chain A suite 54, showed a similar crankshaft motion between test and reference. This suite's probable answer was **5n**, one of the conformer classes not clustered by MINT-AGE (since there were only 4 examples in the HDP33 training set), so a prediction match was impossible in this case. This suggests that this collection of suites instead represents a difficult-to-interpret modeling error idiosyncratic to these ribosome structures, especially 7st2.

Several test suites showed high similarity to the conformations of their respective answer suites, including 2qbz chain X suites 102 and 155, 5b2o B 18, 7st2 3 116 and 564, 6b3k R 58, and 8eiu A 13. In all of the listed cases, the answer suite had a low suiteness score less than 0.1. *Suiteness* is a score between 0 and 1 generated by `phenix.suitename` that reflects the distance between an observed suite conformation and the canonical conformation. The low suiteness scores of these answer suites indicate that the observed conformations are on the border between recognized conformations or are otherwise unusual. We expect that an expanded training set would improve prediction quality around conformer class boundaries. We found that 26% of the non-**1a/1c** prediction mismatches and 18% of the successful prediction matches had an answer suite with suiteness less than 0.1. Therefore, target conformers near class boundaries are somewhat more difficult but are not a systematic impediment to RNAprecis prediction.

Notably absent from the list of prediction mismatches are suite conformers with small MINTAGE clusters. Only one such test suite, namely 2qbz chain X suite 30, has a probable answer whose largest MINT-AGE cluster contains fewer than 5 members, **4g** with 4 members. The prediction match list contains 9 such suites. These suites represent 3% of the mismatches and 5% of the matches, indicating that MINT-AGE successfully produces clusters for low-population conformers, and that MINT-AGE cluster size is not a limiting factor in RNAprecis predictions.

#### 4.5 Other methods

Updates to the suite conformer classes based on new data and MINT-AGE clustering will improve any methods that draw on the consensus conformer definitions. The emphasis in programs similar to ours has shifted. The excellent Erraser [7] appears unavailable at time of writing. Another long existing tool is BrickworX [44], which, at low resolution, often

produces partial structures scattered over multiple symmetry mates and thus appears best suited for the construction of new structures rather than error correction. There are now more aids for getting approximate structures (that is, without consideration of details such as pucker or crankshaft) from as few as a single atom per residue, such as Arena [45]. There have been major developments in prediction, and in the related big calculations of approximate structures from 5–10Å resolution cryoEM data, such as Ribosolve [46]. And there are exciting new cryo-EM structures of ribosomes at very high resolution [34] and of ornate RNA-only structures at about 3Å resolution [47]. But for getting the details right at relatively high resolution, our closest comparison is RCrane [48], which is available and widely used in Coot [49]. It provides structure building of RNA from the phosphate and base and uses the Pperp and backbone conformers from our joint paper [12]. It does an excellent job. However, all these programs are primarily model building tools. In contrast, RNAprecis is focused on model correction and on particularly difficult outliers that persist after refinement. The LDT test set demonstrates that solvable modeling/refinement errors still occur in recent RNA structures. In this sense, the methods are complementary.

## 5 Outlook

Ultimately, we intend to develop RNAprecis predictions into an RNA structure validation and recommendation tool for use in Phenix [10], MolProbity, and similar platforms. Towards this goal, we performed preliminary refinement experiments to replace incorrect conformations in our test-case lower resolution structures with RNAprecis predictions. In most cases our corrections reverted to the incorrect conformation during refinement (see Fig AF in S1 Text) or caused other issues during refinement. That should not be surprising, as a low-resolution map does not contain enough information to determine all atom positions in the RNA backbone. At minimum it requires additional refinement restraints on backbone dihedral angles to take advantage of the extra information from our new set of coordinates based only on the few features really seen at low resolution. In our future plans for RNAprecis, we anticipate building tighter integration with the refinement process, with an automated procedure to create appropriate dihedral restraints to pass to refinement.

Challenges in RNAprecis prediction have identified additional targeted structural tools to be developed alongside RNAprecis. The **1a/1c** ambiguity may yet be resolvable with the RNAprecis parameters such as the  $d_3$  distance. But such difficult cases might also be addressed with paired, comparative local refinements, once RNAprecis identifies a likely site. Automatically identifying suites where the critical RNAprecis atoms are not well-aligned with the map will be a more challenging problem, but would benefit both model building and RNAprecis predictions, especially avoiding pucker mismatches. RNAprecis predictions may be enhanced by employing a more suitable distribution model for clusters. Firstly, a larger training set could make it possible to forego the regularization of cluster covariances, thereby improving the probabilistic interpretability of posterior probabilities as probabilities for a certain cluster. Secondly, the Gaussian likelihood model used here for simplicity may not provide the best fit to the data. It may be worthwhile to explore more general probability distribution models to better capture the data distribution in low-detail coordinates.

MINT-AGE clustering provides an opportunity to revisit and expand previous suite definitions. Suite conformers such as **0a** and **7p** which contain multiple subclusters with distinct enough conformations to merit their own suite identities have been discussed above. Additionally, there may exist additional valid suite conformers that were too rare to be detected at all at the time of the previous work. In the current work, unrecognized suite conformers (!!s) were intentionally excluded from training. MINT-AGE clustering can be performed on the !!s as well, to identify any new clusters of sufficient distribution and consistency to merit recognition as suite conformers, though human inspection will be required to distinguish between real conformations and decoys derived from common modeling errors.

## Supporting information

**S1 Text. Supplement to RNAprecis: Prediction of full-detail RNA conformation from the experimentally best-observed sparse parameters.**

(PDF)

## Author contributions

- Conceptualization:** Henrik Wiechers, Christopher J. Williams, Benjamin Eltzner, Michael G. Prisant, Ezra Miller, Kanti V. Mardia, Jane S. Richardson, Stephan F. Huckemann.
- Data curation:** Christopher J. Williams, Vincent B. Chen, Jane S. Richardson.
- Formal analysis:** Henrik Wiechers, Christopher J. Williams, Benjamin Eltzner, Jane S. Richardson.
- Funding acquisition:** Jane S. Richardson, Stephan F. Huckemann.
- Investigation:** Henrik Wiechers, Christopher J. Williams, Benjamin Eltzner, Jane S. Richardson.
- Methodology:** Henrik Wiechers, Christopher J. Williams, Benjamin Eltzner, Franziska Hoppe, Michael G. Prisant, Vincent B. Chen, Ezra Miller, Jane S. Richardson, Stephan F. Huckemann.
- Project administration:** Jane S. Richardson, Stephan F. Huckemann.
- Resources:** Stephan F. Huckemann.
- Software:** Henrik Wiechers, Benjamin Eltzner.
- Supervision:** Benjamin Eltzner, Kanti V. Mardia, Jane S. Richardson, Stephan F. Huckemann.
- Validation:** Christopher J. Williams, Jane S. Richardson.
- Visualization:** Henrik Wiechers, Christopher J. Williams, Benjamin Eltzner, Franziska Hoppe, Jane S. Richardson.
- Writing – original draft:** Henrik Wiechers, Christopher J. Williams, Benjamin Eltzner, Jane S. Richardson, Stephan F. Huckemann.
- Writing – review & editing:** Christopher J. Williams, Benjamin Eltzner, Michael G. Prisant, Vincent B. Chen, Jane S. Richardson, Stephan F. Huckemann.

## References

1. Steitz TA, Moore PB. RNA, the first macromolecular catalyst: the ribosome is a ribozyme. *Trends Biochem Sci.* 2003;28(8):411–8. [https://doi.org/10.1016/S0968-0004\(03\)00169-5](https://doi.org/10.1016/S0968-0004(03)00169-5) PMID: [12932729](https://pubmed.ncbi.nlm.nih.gov/12932729/)
2. Davis IW, Leaver-Fay A, Chen VB, Block JN, Kapral GJ, Wang X, et al. MolProbity: all-atom contacts and structure validation for proteins and nucleic acids. *Nucleic Acids Res.* 2007;35(Web Server issue):W375–83. <https://doi.org/10.1093/nar/gkm216> PMID: [17452350](https://pubmed.ncbi.nlm.nih.gov/17452350/)
3. Chen VB, Arendall WB 3rd, Headd JJ, Keedy DA, Immormino RM, Kapral GJ, et al. MolProbity: all-atom structure validation for macromolecular crystallography. *Acta Crystallogr D Biol Crystallogr.* 2010;66(Pt 1):12–21. <https://doi.org/10.1107/S0907444909042073> PMID: [20057044](https://pubmed.ncbi.nlm.nih.gov/20057044/)
4. Jain S, Richardson DC, Richardson JS. Computational Methods for RNA Structure Validation and Improvement. *Methods Enzymol.* 2015;558:181–212. <https://doi.org/10.1016/bs.mie.2015.01.007> PMID: [26068742](https://pubmed.ncbi.nlm.nih.gov/26068742/)
5. Keating KS, Pyle AM. Semiautomated model building for RNA crystallography using a directed rotameric approach. *Proc Natl Acad Sci U S A.* 2010;107(18):8177–82. <https://doi.org/10.1073/pnas.0911888107> PMID: [20404211](https://pubmed.ncbi.nlm.nih.gov/20404211/)
6. Read RJ, Adams PD, Arendall WB 3rd, Brunger AT, Emsley P, Joosten RP, et al. A new generation of crystallographic validation tools for the protein data bank. *Structure.* 2011;19(10):1395–412. <https://doi.org/10.1016/j.str.2011.08.006> PMID: [22000512](https://pubmed.ncbi.nlm.nih.gov/22000512/)
7. Chou F-C, Sripakdeevong P, Dibrov SM, Hermann T, Das R. Correcting pervasive errors in RNA crystallography through enumerative structure prediction. *Nat Methods.* 2013;10(1):74–6. <https://doi.org/10.1038/nmeth.2262> PMID: [23202432](https://pubmed.ncbi.nlm.nih.gov/23202432/)
8. Kapral GJ, Jain S, Noeske J, Doudna JA, Richardson DC, Richardson JS. New tools provide a second look at HDV ribozyme structure, dynamics and cleavage. *Nucleic Acids Res.* 2014;42(20):12833–46. <https://doi.org/10.1093/nar/gku992> PMID: [25326328](https://pubmed.ncbi.nlm.nih.gov/25326328/)
9. Adams PD, Afonine PV, Bunkóczi G, Chen VB, Davis IW, Echols N, et al. PHENIX: a comprehensive Python-based system for macromolecular structure solution. *Acta Crystallogr D Biol Crystallogr.* 2010;66(Pt 2):213–21. <https://doi.org/10.1107/S0907444909052925> PMID: [20124702](https://pubmed.ncbi.nlm.nih.gov/20124702/)
10. Liebschner D, Afonine PV, Baker ML, Bunkóczi G, Chen VB, Croll TI, et al. Macromolecular structure determination using X-rays, neutrons and electrons: recent developments in Phenix. *Acta Crystallogr D Struct Biol.* 2019;75(Pt 10):861–77. <https://doi.org/10.1107/S2059798319011471> PMID: [31588918](https://pubmed.ncbi.nlm.nih.gov/31588918/)
11. Murray LJW, Arendall WB 3rd, Richardson DC, Richardson JS. RNA backbone is rotameric. *Proc Natl Acad Sci U S A.* 2003;100(24):13904–9. <https://doi.org/10.1073/pnas.1835769100> PMID: [14612579](https://pubmed.ncbi.nlm.nih.gov/14612579/)

12. Richardson JS, Schneider B, Murray LW, Kapral GJ, Immormino RM, Headd JJ, et al. RNA backbone: consensus all-angle conformers and modular string nomenclature (an RNA Ontology Consortium contribution). *RNA*. 2008;14(3):465–81. <https://doi.org/10.1261/rna.657708> PMID: [18192612](https://pubmed.ncbi.nlm.nih.gov/18192612/)
13. Dryden IL, Mardia KV. *Statistical Shape Analysis*. 2nd ed. Chichester: Wiley; 2016.
14. Mardia KV, Eltzner B, Huckemann SF. *Constrained Shape Analysis with Applications to RNA Structure*. arXiv preprint arXiv:250216270. 2025.
15. Parkinson G, Vojtechovsky J, Clowney L, Brünger AT, Berman HM. New parameters for the refinement of nucleic acid-containing structures. *Acta Crystallogr D Biol Crystallogr*. 1996;52(Pt 1):57–64. <https://doi.org/10.1107/S0907444995011115> PMID: [15299726](https://pubmed.ncbi.nlm.nih.gov/15299726/)
16. Florek K, Łukasiewicz J, Perkal J, Steinhaus H, Zubrzycki S. Sur la liaison et la division des points d'un ensemble fini. *Colloq Math*. 1951;2(3–4):282–5. <https://doi.org/10.4064/cm-2-3-4-282-285>
17. Langfelder P, Zhang B, Horvath S. Defining clusters from a hierarchical cluster tree: the Dynamic Tree Cut package for R. *Bioinformatics*. 2008;24(5):719–20. <https://doi.org/10.1093/bioinformatics/btm563> PMID: [18024473](https://pubmed.ncbi.nlm.nih.gov/18024473/)
18. Obulkasim A, Meijer GA, van de Wiel MA. Semi-supervised adaptive-height snipping of the hierarchical clustering tree. *BMC Bioinformatics*. 2015;16(1):15. <https://doi.org/10.1186/s12859-014-0448-1> PMID: [25592847](https://pubmed.ncbi.nlm.nih.gov/25592847/)
19. Mardia KV, Wiechers H, Eltzner B, Huckemann SF. Principal component analysis and clustering on manifolds. *J Multivar Anal*. 2021;188:104862. <https://doi.org/10.1016/j.jmva.2021.104862>
20. Eltzner B, Huckemann S, Mardia KV. Torus principal component analysis with applications to RNA structure. *Ann Appl Stat*. 2018;12(2):1332–59. <https://doi.org/10.1214/17-AOAS1115>
21. Jung S, Dryden IL, Marron JS. Analysis of principal nested spheres. *Biometrika*. 2012;99(3):551–68. <https://doi.org/10.1093/biomet/ass022> PMID: [23843669](https://pubmed.ncbi.nlm.nih.gov/23843669/)
22. Oliveira M, Crujeiras RM, Rodríguez-Casal A. CircSiZer: an exploratory tool for circular data. *Environ Ecol Stat*. 2013;21(1):143–59. <https://doi.org/10.1007/s10651-013-0249-0>
23. Huckemann S, Kim K-R, Munk A, Rehfeldt F, Sommerfeld M, Weickert J, et al. The circular SiZer, inferred persistence of shape parameters and application to early stem cell differentiation. *Bernoulli*. 2016;22(4). <https://doi.org/10.3150/15-bej722>
24. Dümbgen L, Walther G. Multiscale inference about a density. *Ann Statist*. 2008;36(4):1758–85. <https://doi.org/10.1214/07-AOS521>
25. Chaudhuri P, Marron JS. Scale space view of curve estimation. *Ann Statist*. 2000;28(2):408–28. <https://doi.org/10.1214/aos/1016218224>
26. Huckemann SF, Eltzner B, Mardia KV. RNA Structure Correction – the Importance of Small Clusters. In: *Geometric Science of Information 2025 proceedings*. Springer; 2025. p. to appear.
27. Bookstein FL. The measurement of biological shape and shape change. vol. 24 of *Lecture Notes in Biomathematics*. 2nd ed. New York: Springer-Verlag; 1978.
28. Goodall CR, Mardia KV. Multivariate Aspects of Shape Theory. *Ann Statist*. 1993;21(2). <https://doi.org/10.1214/aos/1176349154>
29. Mardia KV, Patrangenaru V. On affine and projective shape data analysis. *Functional and Spatial Data Analysis, Proceedings of the 20th LASR Workshop* (Eds: KV Mardia and RG Aykroyd). 2001; p. 39–45.
30. Fréchet M. Les éléments aléatoires de nature quelconque dans un espace distancié. *Ann Inst Henri Poincaré*. 1948;10(4):215–310.
31. Wiechers H, Eltzner B, Mardia KV, Huckemann SF. Learning torus PCA-based classification for multiscale RNA correction with application to SARS-CoV-2. *J R Stat Soc Ser C Appl Stat*. 2023;72(2):271–93. <https://doi.org/10.1093/jrsssc/qlad004>
32. Williams C, Richardson J. High-quality RNA residues: RNA2023; 2023. Available from: <https://doi.org/10.5281/zenodo.8103014>
33. Leontis NB, Zirbel CL. Nonredundant 3D structure datasets for RNA knowledge extraction and benchmarking. In: *Leontis NB, Westhof E, editors. RNA 3D Structure Analysis and Prediction*. Springer; 2012. p. 281–98.
34. Fromm SA, O'Connor KM, Purdy M, Bhatt PR, Loughran G, Atkins JF, et al. The translating bacterial ribosome at 1.55 Å resolution generated by cryo-EM imaging services. *Nat Commun*. 2023;14(1):1095. <https://doi.org/10.1038/s41467-023-36742-3> PMID: [36841832](https://pubmed.ncbi.nlm.nih.gov/36841832/)
35. Faille A, Dent KC, Pellegrino S, Jaako P, Warren AJ. The chemical landscape of the human ribosome at 1.67 Å resolution. *bioRxiv*. 2023. <https://doi.org/10.1101/2023.02.28.530191>
36. Williams CJ, Richardson DC, Richardson JS. The importance of residue-level filtering and the Top2018 best-parts dataset of high-quality protein residues. *Protein Sci*. 2022;31(1):290–300. <https://doi.org/10.1002/pro.4239> PMID: [34779043](https://pubmed.ncbi.nlm.nih.gov/34779043/)
37. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, et al. The Protein Data Bank. *Nucleic Acids Res*. 2000;28(1):235–42. <https://doi.org/10.1093/nar/28.1.235> PMID: [10592235](https://pubmed.ncbi.nlm.nih.gov/10592235/)
38. Mardia KV, Jupp PE. *Directional Statistics*. New York: Wiley; 2000.
39. Arnaudon M, Miclo L. Means in complete manifolds: uniqueness and approximation. *ESAIM: PS*. 2014;18:185–206. <https://doi.org/10.1051/ps/2013033>
40. Mahalanobis PC. On the generalized distance in statistics. *Sankhyā A*. 1948;10:1–7.
41. Huang L, Wang J, Lilley DMJ. A critical base pair in k-turns determines the conformational class adopted, and correlates with biological function. *Nucleic Acids Res*. 2016;44(11):5390–8. <https://doi.org/10.1093/nar/gkw201> PMID: [27016741](https://pubmed.ncbi.nlm.nih.gov/27016741/)

42. Klein DJ, Schmeing TM, Moore PB, Steitz TA. The kink-turn: a new RNA secondary structure motif. *EMBO J.* 2001;20(15):4214–21. <https://doi.org/10.1093/emboj/20.15.4214> PMID: [11483524](https://pubmed.ncbi.nlm.nih.gov/11483524/)
43. Leontis NB, Westhof E. A common motif organizes the structure of multi-helix loops in 16 S and 23 S ribosomal RNAs. *J Mol Biol.* 1998;283(3):571–83. <https://doi.org/10.1006/jmbi.1998.2106> PMID: [9784367](https://pubmed.ncbi.nlm.nih.gov/9784367/)
44. Chojnowski G, Waleń T, Piątkowski P, Potrzebowski W, Bujnicki JM. Brickworx builds recurrent RNA and DNA structural motifs into medium- and low-resolution electron-density maps. *Acta Crystallogr D Biol Crystallogr.* 2015;71(Pt 3):697–705. <https://doi.org/10.1107/S1399004715000383> PMID: [25760616](https://pubmed.ncbi.nlm.nih.gov/25760616/)
45. Perry ZR, Pyle AM, Zhang C. Arena: Rapid and Accurate Reconstruction of Full Atomic RNA Structures From Coarse-grained Models. *J Mol Biol.* 2023;435(18):168210. <https://doi.org/10.1016/j.jmb.2023.168210> PMID: [37479079](https://pubmed.ncbi.nlm.nih.gov/37479079/)
46. Kappel K, Zhang K, Su Z, Watkins AM, Kladwang W, Li S, et al. Accelerated cryo-EM-guided determination of three-dimensional RNA-only structures. *Nat Methods.* 2020;17(7):699–707. <https://doi.org/10.1038/s41592-020-0878-9> PMID: [32616928](https://pubmed.ncbi.nlm.nih.gov/32616928/)
47. Kretsch RC, Wu Y, Shabalina SA, Lee H, Nye G, Koonin EV, et al. Naturally ornate RNA-only complexes revealed by cryo-EM. *Nature.* 2025;643(8073):1135–42. <https://doi.org/10.1038/s41586-025-09073-0> PMID: [40328315](https://pubmed.ncbi.nlm.nih.gov/40328315/)
48. Keating KS, Pyle AM. RCrane: semi-automated RNA model building. *Acta Crystallogr D Biol Crystallogr.* 2012;68(Pt 8):985–95. <https://doi.org/10.1107/S0907444912018549> PMID: [22868764](https://pubmed.ncbi.nlm.nih.gov/22868764/)
49. Emsley P, Lohkamp B, Scott WG, Cowtan K. Features and development of Coot. *Acta Crystallogr D Biol Crystallogr.* 2010;66(Pt 4):486–501. <https://doi.org/10.1107/S0907444910007493> PMID: [20383002](https://pubmed.ncbi.nlm.nih.gov/20383002/)