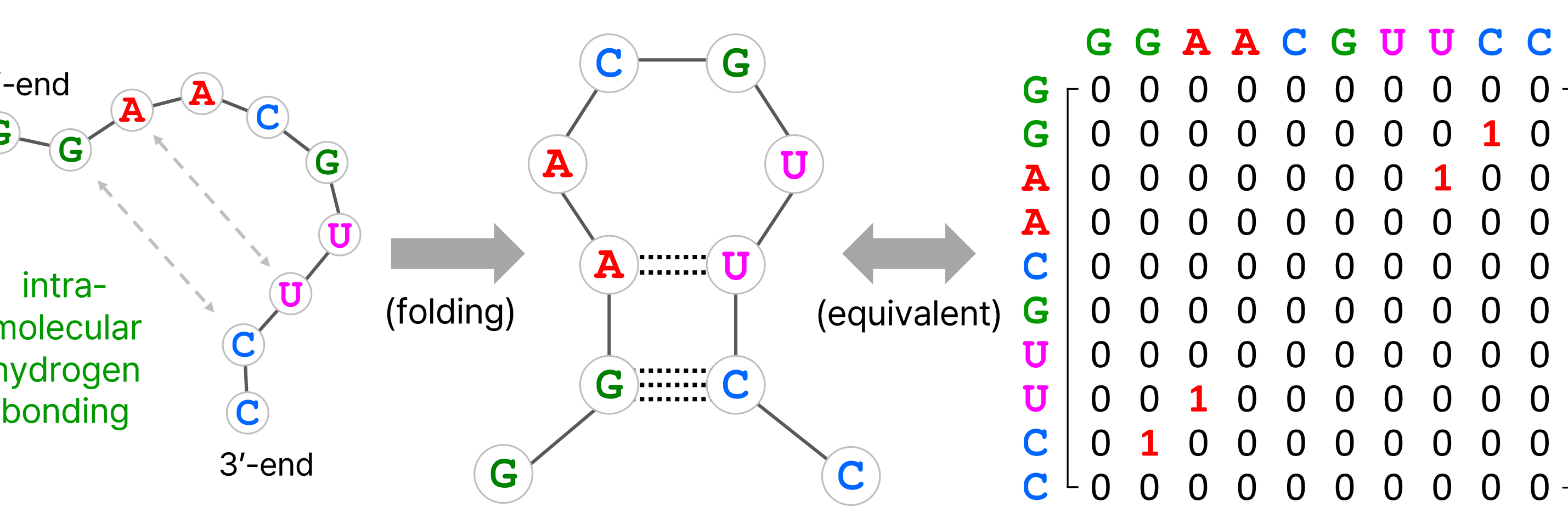


RNA Secondary Structure



RNA strand (5'-end to 3'-end) undergoes (folding) to form Secondary structure, which is (equivalent) to a Binary matrix.

Q. Do RNA structures matter? If yes, why?

A. Yes, as **structure is the main determinant of the function of RNAs**.
→ Knowing RNA structure is crucial for advances in biotechnology.

Problem Setting

Primary structure (as a 1D sequence of nucleotides) $\mathbf{x} = (x_1, \dots, x_L)$ is input to an ML Model to produce a Score matrix $\hat{\mathbf{Y}} \in \mathbb{R}^{L \times L}$ (as a 2D matrix).

Train such a model with a collection of experimentally validated samples

$$\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{Y}^{(i)})\}_{i=1}^m$$

measured via X-ray crystallography / nuclear magnetic resonance

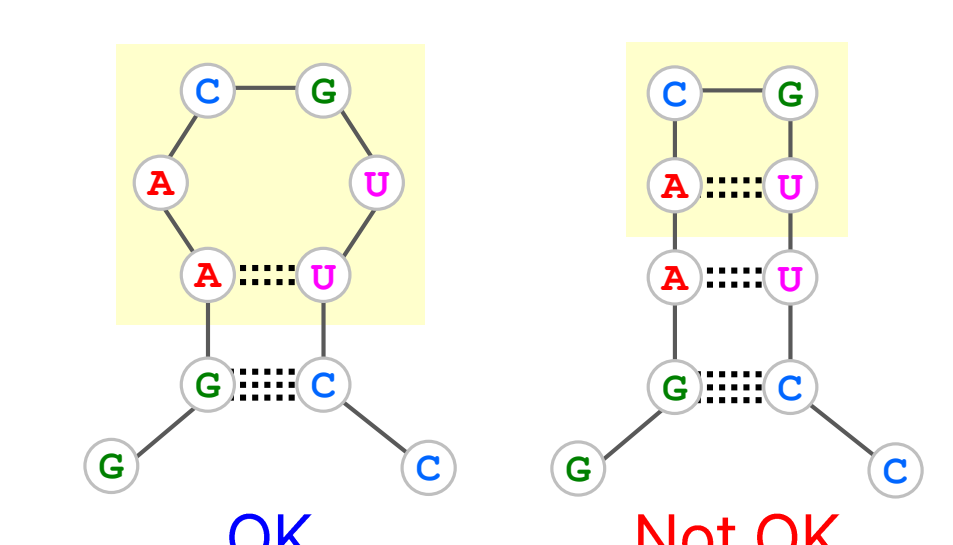
Considering the **natural constraints** governed by *physical laws*:

(C1) Binary & symmetric
 $Y_{ij} \in \{0, 1\}$ for $\forall i, j$ and $\mathbf{Y} = \mathbf{Y}^T$,

(C2) Watson-Crick & Wobble base pairs only
 $Y_{ij} = 0$ if $x_i x_j \notin \mathcal{B} := \{\text{AU, UA, GC, CG, GU, UG}\}$

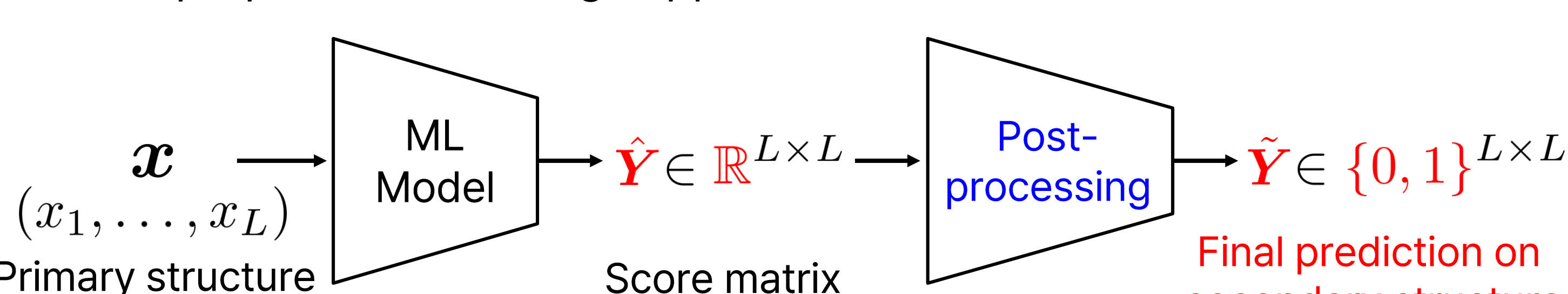
(C3) No sharp loops
 $Y_{ij} = 0$ if $|i - j| < 4$

(C4) No overlap of pairs



Prior Works & Challenge

E2efold proposes a two-stage approach:



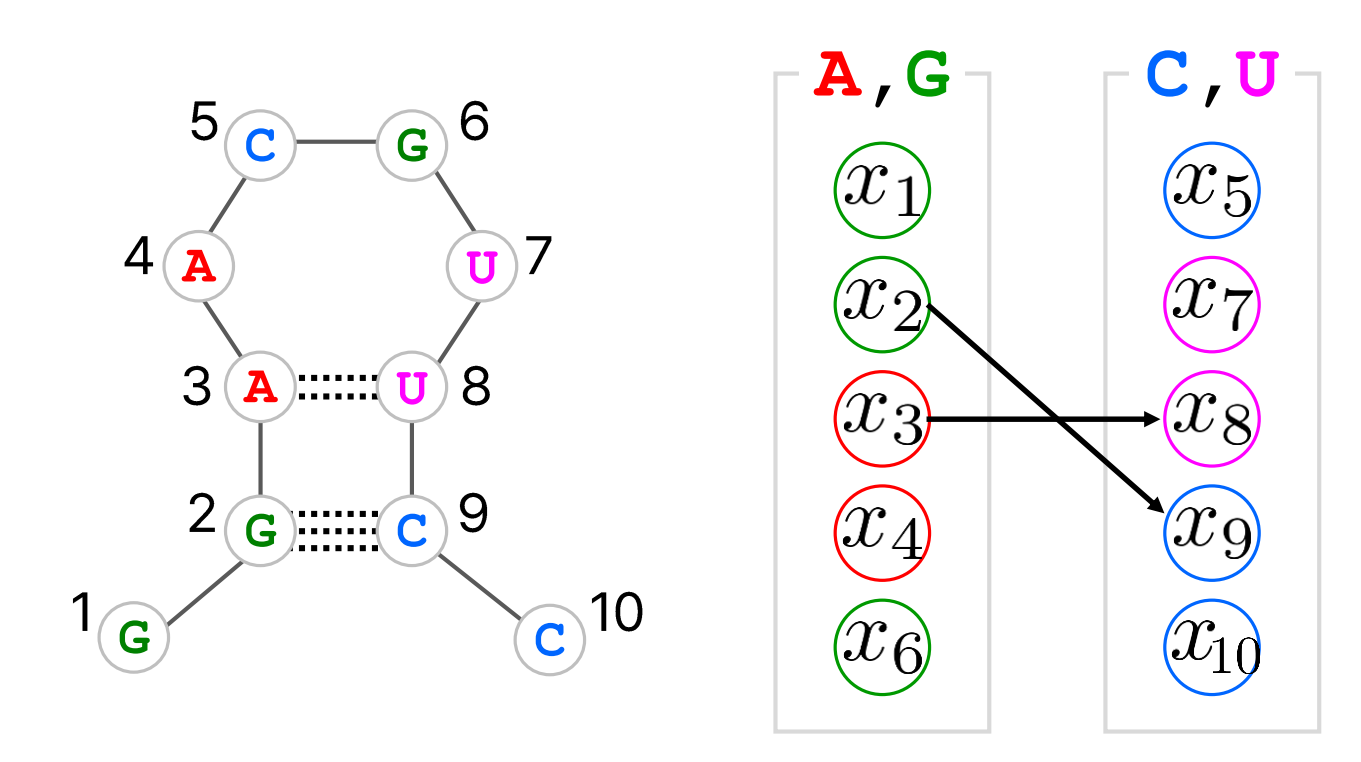
where post-processing for the constraints is presented as an optimization:

$$\begin{aligned} \max_{\tilde{\mathbf{Y}}} \quad & \langle \hat{\mathbf{Y}} - s, \tilde{\mathbf{Y}} \rangle \\ \text{s.t.} \quad & \tilde{Y}_{ij} \in \{0, 1\}, \quad \tilde{\mathbf{Y}} = \tilde{\mathbf{Y}}^T, \quad \tilde{\mathbf{Y}} \mathbf{1} \leq \mathbf{1}, \\ & \tilde{Y}_{ij} = 0 \quad \text{if } x_i x_j \notin \mathcal{B} \text{ or } |i - j| < 4. \end{aligned}$$

Predictions made by **SOTA algorithms** still violate (C2) – (C4)
 such as CNNfold, E2efold, Ufold, REDfold

Contributions

Propose a mathematically equivalent post-process optimization based on the **assignment problem**



$$\begin{aligned} \min_{\mathbf{Z}} \quad & \langle \mathbf{C}, \mathbf{Z} \rangle \quad \text{assignment matrix} \\ \text{s.t.} \quad & \mathbf{Z}_{ij} \in \{0, 1\}, \\ & \mathbf{Z} \mathbf{1} = \mathbf{1}, \quad \mathbf{Z}^T \mathbf{1} = \mathbf{1}. \end{aligned}$$

Theorem 1 (Equivalence)

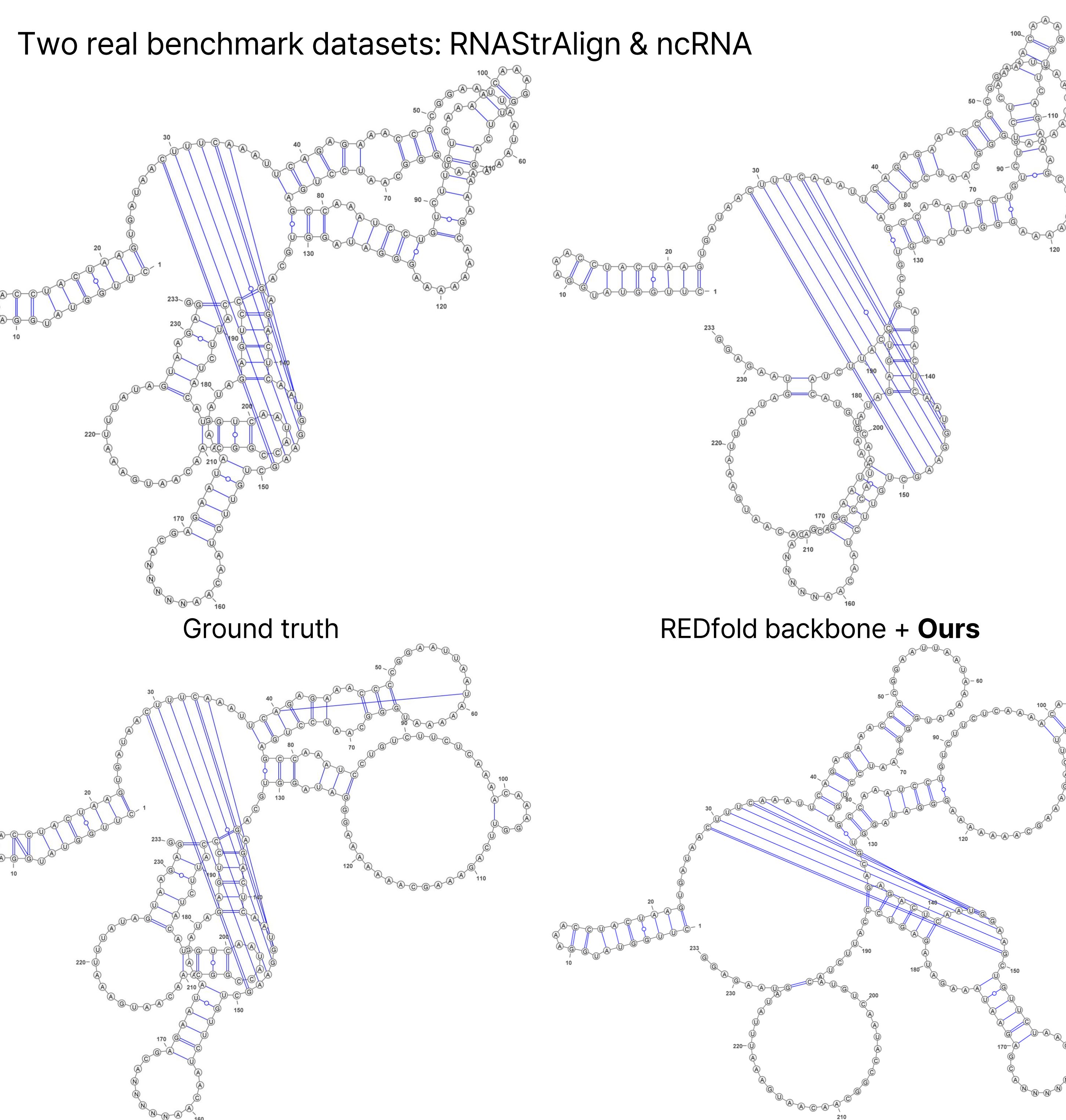
Any optimal solution of each optimization can be reduced to an optimal solution of the other problem.

→ Employ off-the-shelf algorithms such as **Hungarian algorithm / Zonker-Volgenant algorithm** s.t. we obtain:

- Model-agnostic add-on to any backbone ML model
- Outputs completely adhering to the fundamental constraints of RNAs
- Improved predictive performance & empirical running time

Results

Two real benchmark datasets: RNAstrAlign & ncRNA



Method	Constraint violation			Prediction performance			Run time
	(C2)	(C3)	(C4)	F1	Recall	Precision	Time (s)
RNAfold	-	-	-	0.606	0.679	0.566	-
RNAstructure	-	-	-	0.599	0.668	0.562	-
CONTRAFold	0%	8.1%	0%	0.626	0.690	0.596	-
SPOT-RNA	57.7%	30.3%	0%	0.647	0.683	0.640	-
MXfold2	-	-	-	0.631	0.687	0.608	-
E2efold + E2E PP	0%	12.8%	58.3%	0.595	0.575	0.631	0.049
E2efold + Blossom	73.0%	21.0%	0.3%	0.489	0.615	0.415	2.212
E2efold + Ours	0%	0%	0%	0.608	0.602	0.622	0.308
REDfold + RED PP	0%	0.5%	0%	0.844	0.849	0.877	0.053
REDfold + Blossom	9.4%	0.3%	0%	0.840	0.873	0.838	0.378
REDfold + Ours	0%	0%	0%	0.847	0.867	0.858	0.005