

# LLM Agent–Driven ncRNA Design via Intrinsic Features and Structure-Guided Feedback

**Work in Progress, Exploratory Analysis**

Preliminary Data — Subject to Discussion and Revision

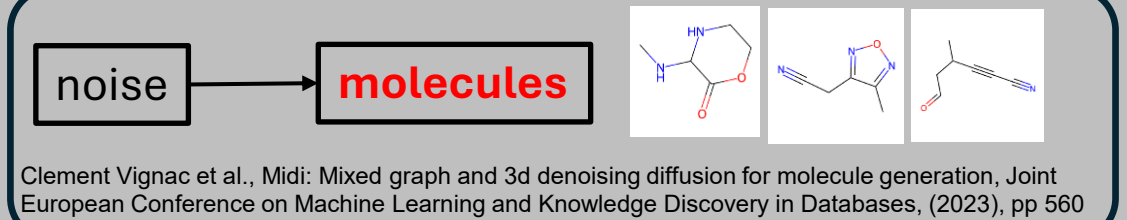
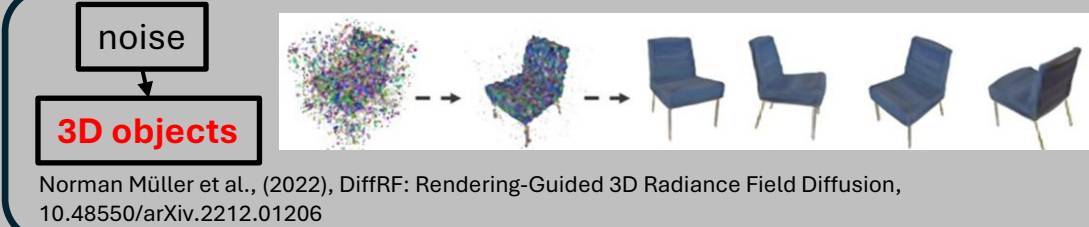
Focuses on RNA ribonucleotide sequences; no DNA or proteins explicitly involved

# Inspiration (from Diffusion Models) and the **Objective**

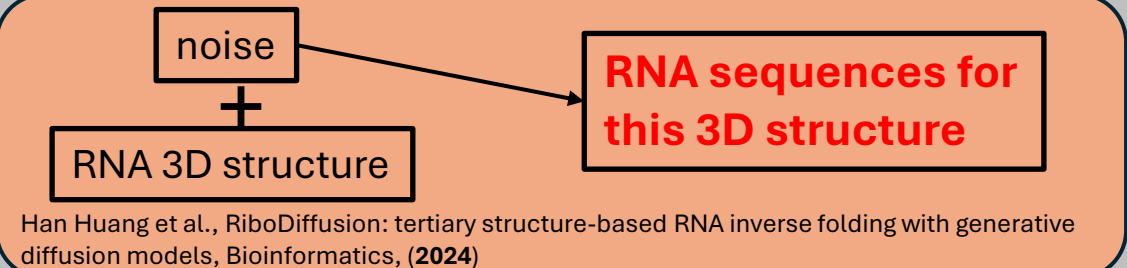
1



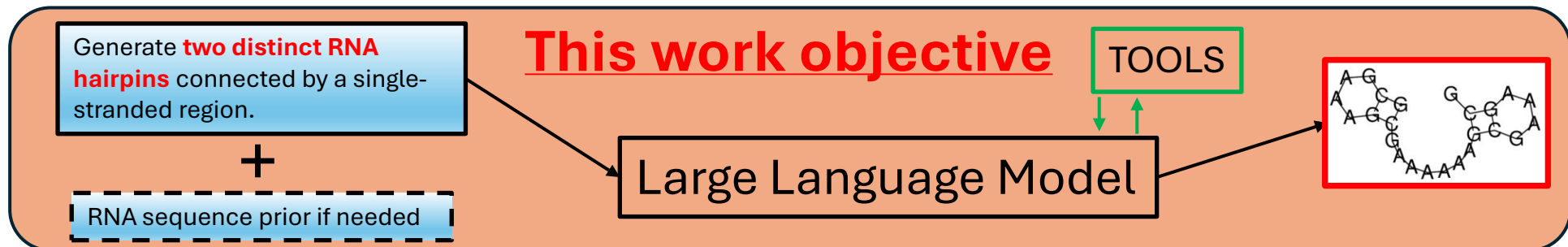
2



3



**objective**



**objective**

# OBJECTIVES

**Develop an LLM agent-driven pipeline that generates and analyzes RNA sequences, under guidance from human-designed prompts. Each pipeline automatically integrates RNA structural feedback, diffusion-based generative model and internet search.**

- **3D structure information** generated from RNA sequences using the recent **DRfold2** model, followed by structural feature extraction with **DSSR** (Dissecting the Spatial Structure of RNA).
- **Structural refinement** is guided by the **conditional diffusion model RiboDiffusion**, that proposes alternative RNA sequences with the same 3D structure as the reference one.

## Literature used:

**LLM:** Anthropic. Claude Sonnet 4 (20250514). <https://www.anthropic.com>. (Mai **2025**)

**LLM ReAct agent:** Shunyu Yao et al., ReAct: Synergizing Reasoning and Acting in Language Models, (**2023**)

**LangGraph (2024):** Low-level orchestration framework for building, managing, and deploying long-running, stateful agents, <https://github.com/langchain-ai/langgraph>

**DRfold2:** Li, Yang Li et al., Ab initio RNA structure prediction with composite language model and denoised end-to-end learning}, Cold Spring Harbor Laboratory, (**2025**)

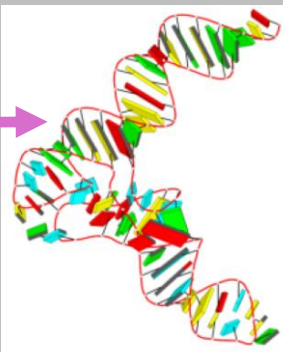
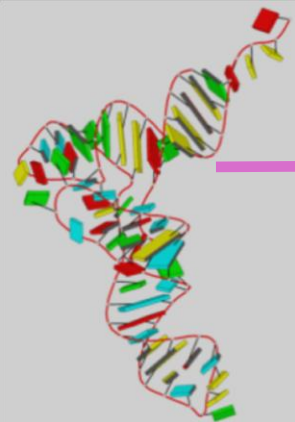
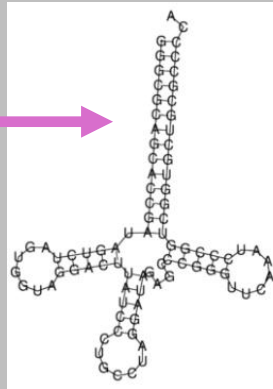
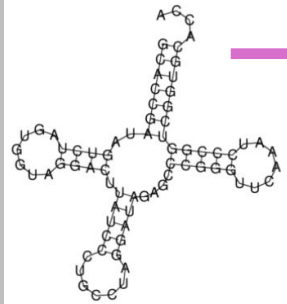
**Ribodiffusion:** Han Huang et al., RiboDiffusion: tertiary structure-based RNA inverse folding with generative diffusion models, Bioinformatics, (**2024**)

**RNA-FM:** Jiayang Chen et al., Interpretable RNA Foundation Model from Unannotated Data for Highly Accurate RNA Structure and Function Predictions, bioRxiv, (**2022**)

**DSSR:** Xiang-Jun Lu, DSSR-enabled innovative schematics of 3D nucleic acid structures with PyMOL, Nucleic Acids Research, (**2020**)

reference sequence      modified sequence

tRNA<sup>Gly</sup>  
Asgard\_group\_archaeon



## LLM ReAct Agent: Modification of a Specified Region of an RNA Molecule

f"Current tRNA sequence (input): {tRNA}\n\n"

"Modify this reference sequence by increasing the length of the acceptor stem to **\*\*twice its original size\*\***. "

"Ensure that **\*\*no more than 20%\*\*** of the nucleotides are changed or added.\n\n"

"Keep the rest of the tRNA unchanged.\n\n"

"Verify the structure and return the modified RNA sequence."

I'll help you modify the tRNA sequence by doubling the acceptor stem length while keeping changes to no more than 20% of nucleotides. First, let me analyze the current structure to identify the acceptor  
Tool Calls: **get\_dotBracket\_length**

'Now let me identify the stems to locate the acceptor stem  
Tool Calls: **get\_stems**

'Let me also get the hairpin loops to better understand the structure  
Tool Calls: **get\_hairpins**

Based on the analysis, I can see this is a standard tRNA structure with:\n- **\*\*Acceptor stem\*\***: **helix\_0** with 7 base pairs (GCACCGA paired with CGUGGCU)\n- **D-arm**: **helix\_1** with 4 base pairs\n- **Anticodon arm**: **helix\_2** with 5 base pairs\n- **TψC arm**: **helix\_3** with 5 base pairs\n\nTo double the acceptor stem length from 7 to 14 base pairs while changing ≤20% of nucleotides (≤15 out of 74), I'll:\n1. Extend the 5' end with 7 complementary base pairs\n2. Extend the 3' end with 7 complementary base pairs\n3. Keep all other structural elements unchanged\n\nHere's the modified sequence: **GGGC...**

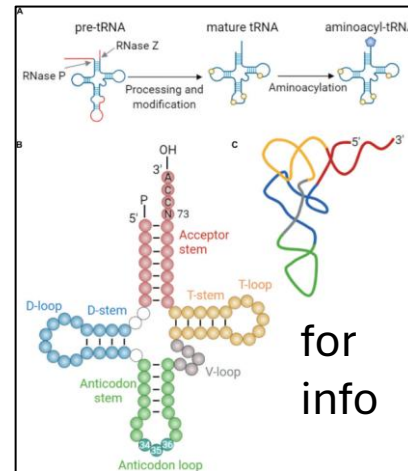


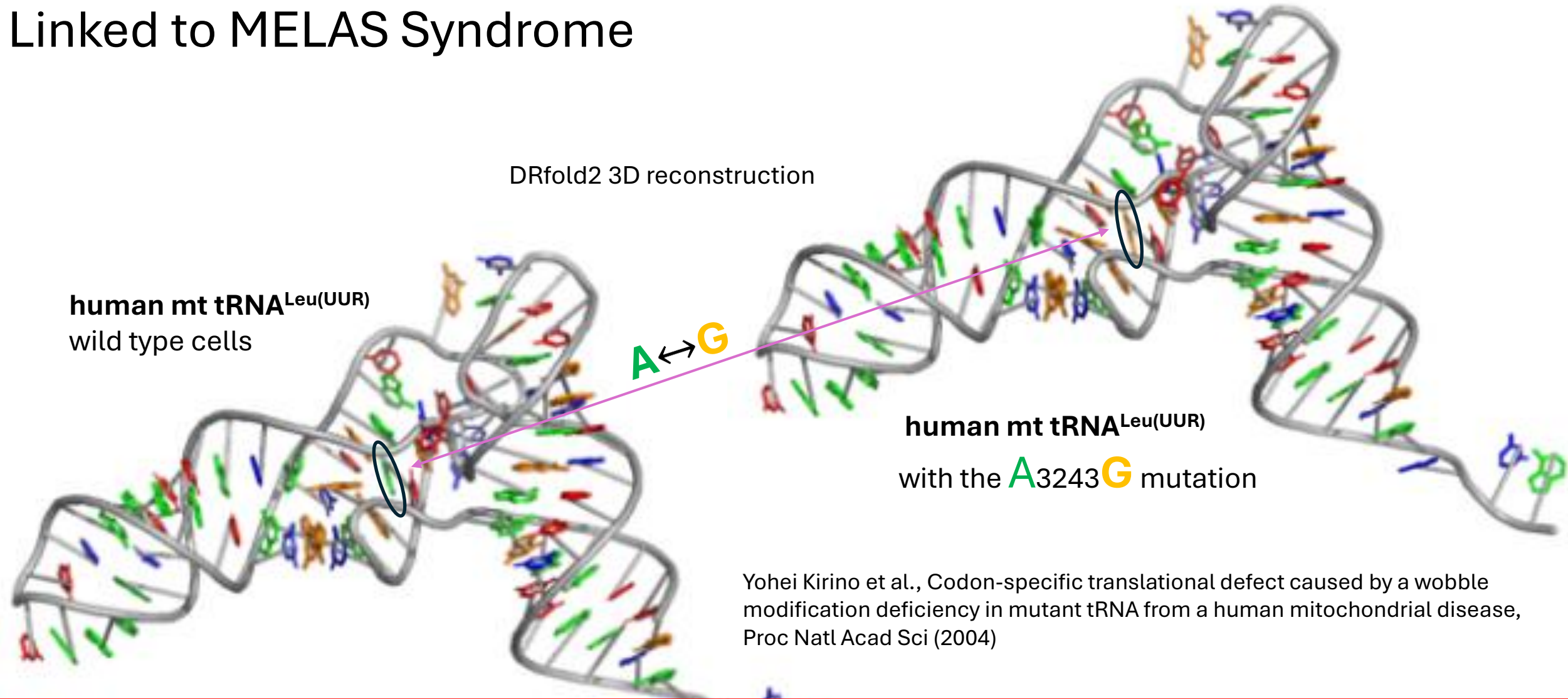
Figure adapted from Frontiers in Molecular Biosciences (Irem Avcar-

This is a slide from the results section.

acceptor stem## Modified tRNA sequence.

`GGGCGCAGCACCGAUAGUCUAGUGGUAGGACUUAUCCCGGCCUAGGAUAGAGCCCGGGUUCAAAUCCCGGUCGGUGCUGCGCCCCA` \*\*The acceptor stem is now twice its original length (14 bp vs 7 bp) while preserving the overall tRNA secondary structure and staying within the 20% nucleotide change constraint.

# Small Mutation, Significant Impact: tRNA Point Mutation (A3243G) Linked to MELAS Syndrome



This is another slide from the results section.



# OUTLINE

## ❖ Dataset

- What is RNA molecule
- RNACentral database

## ❖ Tools description

- RNA alignment Tool
- RNA 3D structure parameters Tool
- The structure-guided feedback Tool

## ❖ LLM **Tool-Agent** Computational Pipeline

### ❖ Results for LLM Tool-Agent

- Speculative simulation of reverse evolution of tRNA<sup>Gly</sup> (6 slides)

## ❖ LLM **ReAct-Agent** Computational Pipeline

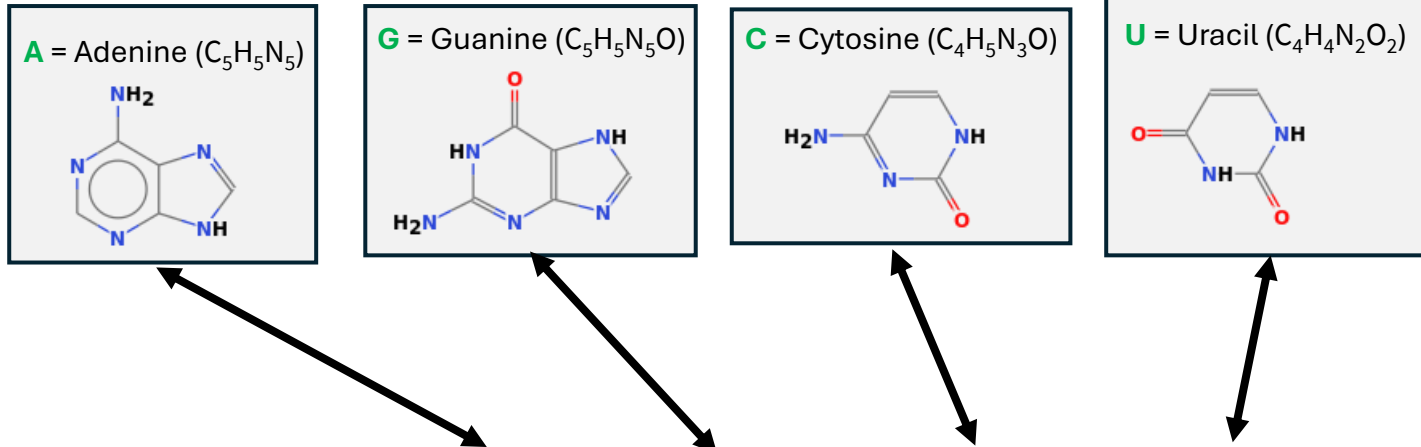
### ❖ Results for LLM ReAct-Agent

- Comparison of two RNAs
- Single point mutation in tRNA associated with MELAS syndrome
- Generate two connected RNA hairpins without RNA template
- Generate two connected RNA hairpins from a given RNA template
- Targeted modification of a specific region of an RNA molecule

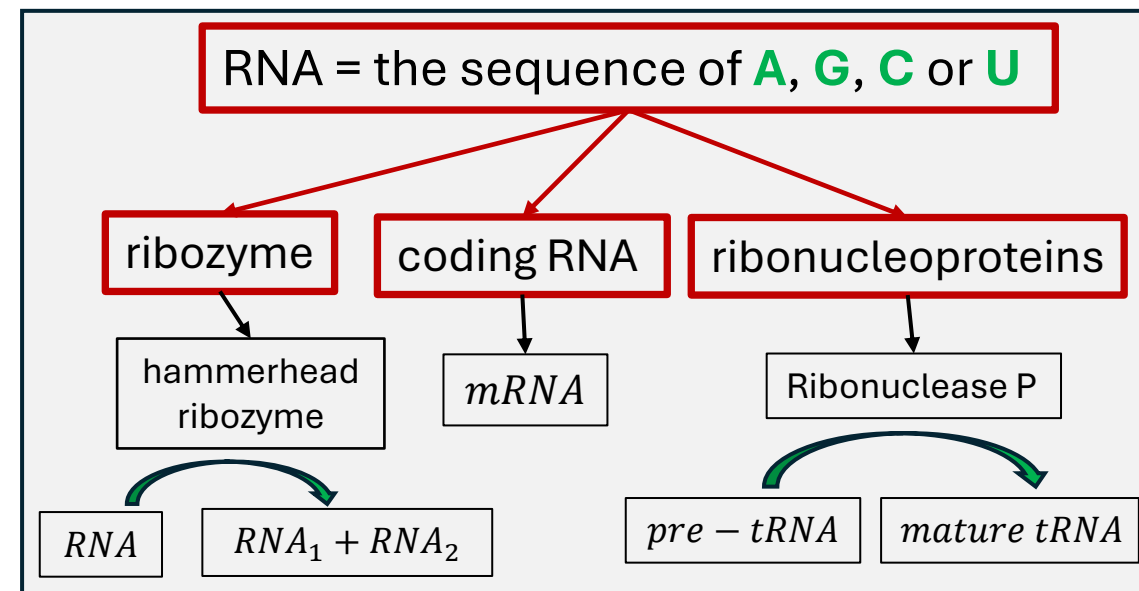
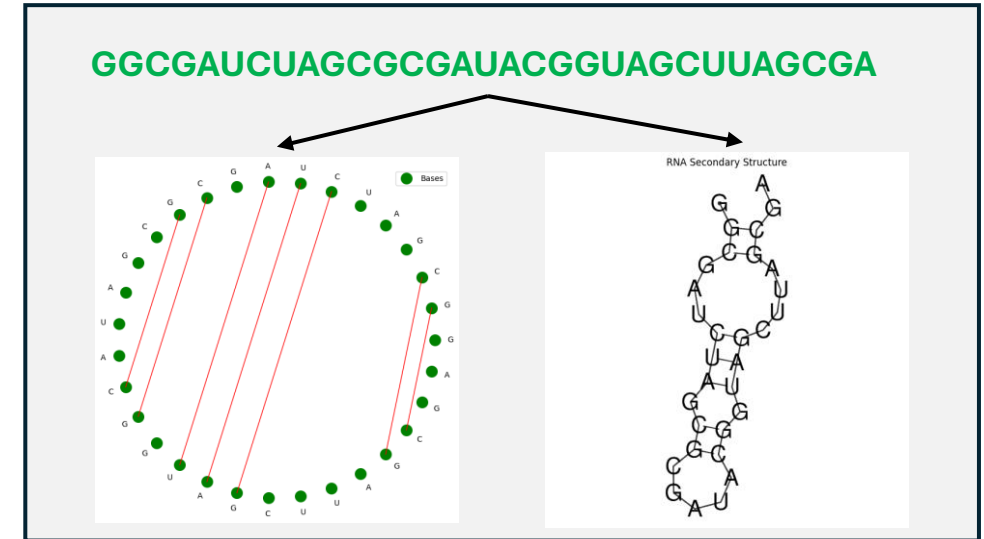
## ❖ Conclusions

## ❖ What's next ?

# What is RNA molecule ?



	Adenine	Guanine	Cytosine	Uracil
<b>O</b>	0	1	1	2
<b>N</b>	5	5	3	2
<b>C</b>	5	5	4	4
<b>H</b>	5	5	5	4
$\Delta_f H_{solid}^0, kJ/mol$	96.9	-183.9	-221	-424.4
$\Delta_c H_{solid}^0, kJ/mol$	-2779.0	-2498.2	-2067	-1721.3
$M_w, g/mol$	135	151	111	112
Hydrogen bonds	2	3	3	2

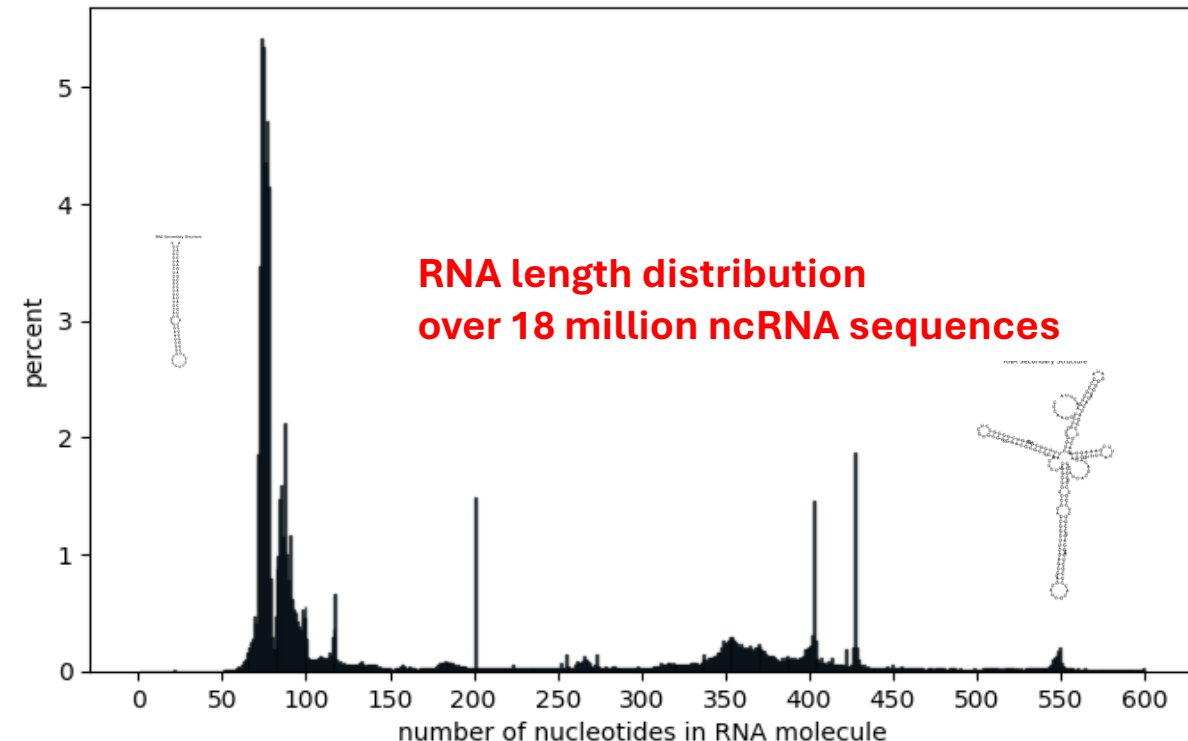
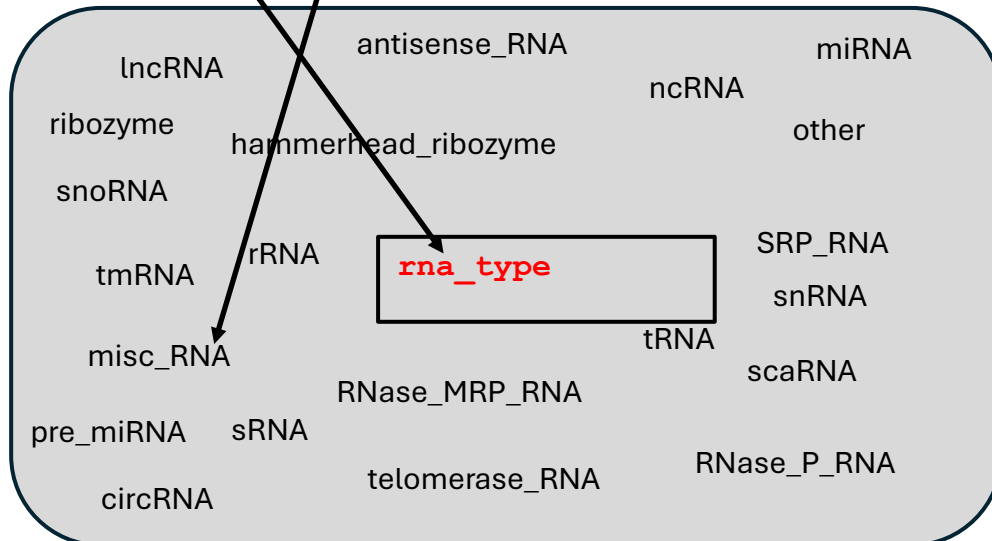
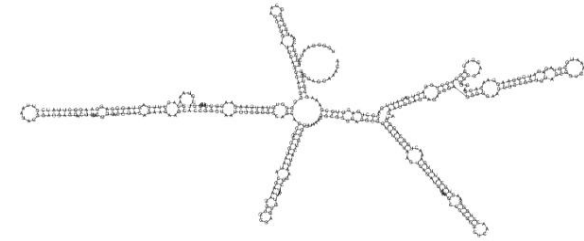


# The dataset: RNACentral database

## RNACentral Browsable API

[https://rnacentral.org/api/v1/rna/?page=3&page\\_size=100](https://rnacentral.org/api/v1/rna/?page=3&page_size=100) gives:

```
{ "url": "http://rnacentral.org/api/v1/rna/URS0002915621",
  "rnacentral_id": "URS0002915621",
  "md5": "fee3fe68dbd91ee898bffd9d4b89b2e9",
  "sequence": "AUGGAUGGUUGAUCAGAGAACGUACAUUUUAUAAAUGGUGUAUGUCAAUUGAUCCACAGUCCCU",
  "length": 64,
  "xrefs": "http://rnacentral.org/api/v1/rna/URS0002915621/xrefs",
  "publications": "http://rnacentral.org/api/v1/rna/URS0002915621/publications",
  "is_active": true,
  "description": "pre_miRNA from 0 species",
  "rna_type": "pre_miRNA",
  "count_distinct_organisms": 4,
  "distinct_databases": [ "Rfam" ] }, ...
```





# OUTLINE

## ❖ Objectives

## ❖ Dataset

- What is RNA molecule
- RNACentral database

## ❖ Tools description

- RNA alignment Tool
- RNA 3D structure parameters Tool
- The structure-guided feedback Tool

## ❖ LLM **Tool-Agent** Computational Pipeline

## ❖ Results for LLM Tool-Agent

- Speculative simulation of reverse evolution of tRNA<sup>Gly</sup> (6 slides)

## ❖ LLM **ReAct-Agent** Computational Pipeline

## ❖ Results for LLM ReAct-Agent

- Comparison of two RNAs
- Single point mutation in tRNA associated with MELAS syndrome
- Generate two connected RNA hairpins without RNA template
- Generate two connected RNA hairpins from a given RNA template
- Targeted modification of a specific region of an RNA molecule

## ❖ Conclusions

## ❖ What's next ?

# RNA alignment Tool:

## RNA Prompt analysis using Large Language Model

The answer is fully correct.

You are an expert in molecular biology. Your task is to analyze two RNA sequences and determine which is more likely to be found in a real biological organism. Justify your answer in one sentence.

Sequence 1: GGGGGGGG

Sequence 2: AAGUCGCGCCGAAAAGGUGUCUCUU

Question: Which sequence is more biologically realistic?

Answer format: One sentence explaining which sequence is more likely and why.

**1. Sequence 2 is more biologically realistic because it has a more diverse nucleotide composition and is not a repetition of a single nucleotide.**

Including additional context would help make the answer clearer and more specific

```
pip install biopython
from Bio.Align import PairwiseAligner
```

```
<class 'Bio.Align.Alignment'>
target      0 GAG--GCG-----GGUG----- 10
            0 .||--|||-----|||----- 25
query       0 AAGUCGCGCCGAAAAGGUGUCUCUU 25
```

The LLM has difficulty understanding this.

```
{
  "target_sequence": "GAGGCGGGUG",
  "query_sequence": "AAGUCGCGCCGAAAAGGUGUCUCUU",
  "aligned_target": "GAG--GCG-----GGUG-----",
  "aligned_query": "AAGUCGCGCCGAAAAGGUGUCUCUU",
  "score": 0.2,
  "identical_nucleotide_counts": {
    "A-A": 1,
    "G-G": 6,
    "C-C": 1,
    "U-U": 1
  }
}
```

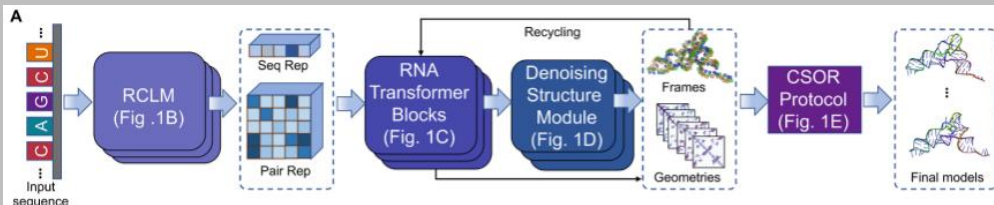
dictionaries **are better understood** by LLMs

USE for  
LLM  
prompt

# RNA 3D structure parameters Tool:

combine DRfold2 and DSSR to extract 3D structural parameters of RNA

DRfold2 enables the prediction of RNA 3D structures (in PDB format) from nucleotide sequences provided in FASTA format



DRfold2

GCGCGCAUACGUGCGCGC

Li, Yang Li et al., Ab initio RNA structure prediction with composite language model and denoised end-to-end learning, Cold Spring Harbor Laboratory, (2025), doi = {10.1101/2025.03.05.641632}

DSSR

API: <http://skmatic.x3dna.org/api>

```
{
  "general_info": {
    "full_RNA_sequence": "GCGCGCAUACGUGCGCGC",
    "length": 18,
    "base_pairs": 7,
    "hydrogen_bonds": 25,
    "dot_bracket": [
      "(((((((.....)))))))))"
    ],
    "splayUnits": [
      "UA"
    ],
    "hairpins": [
      "AUACGU"
    ],
    "stacks": [
      "AU",
      "ACGU"
    ]
  },
  "helices_info": {
    "full_RNA_sequence": "GCGCGCAUACGUGCGCGC",
    "helix_0": {
      "base_pairs": 7,
      "strand_1": "GCGCGCA",
      "strand_2": "CGCGCGU",
      "helix_form": "AAAAAA"
    }
  }
}
```

Xiang-Jun Lu, DSSR-enabled innovative schematics of 3D nucleic acid structures with PyMOL, Nucleic Acids Research, Vol 48, number 13, p. e74-e74, (2020), doi = {10.1093/nar/gkaa426}

Input

Output for LLM input

From RNA sequence to a parametric description of its 3D structure

# The Structure-Guided Feedback Tool: the inverse folding problem

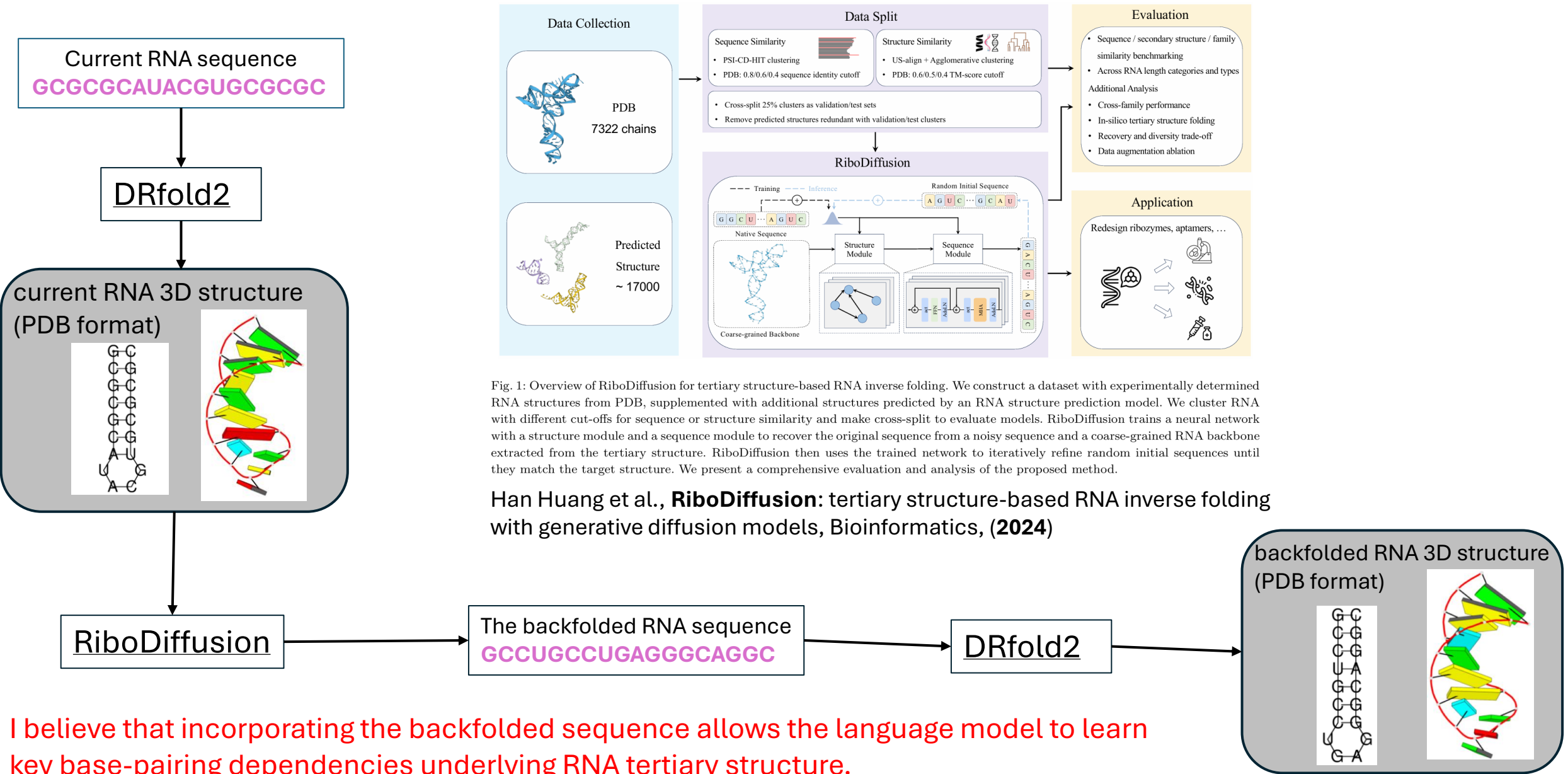


Fig. 1: Overview of RiboDiffusion for tertiary structure-based RNA inverse folding. We construct a dataset with experimentally determined RNA structures from PDB, supplemented with additional structures predicted by an RNA structure prediction model. We cluster RNA with different cut-offs for sequence or structure similarity and make cross-split to evaluate models. RiboDiffusion trains a neural network with a structure module and a sequence module to recover the original sequence from a noisy sequence and a coarse-grained RNA backbone extracted from the tertiary structure. RiboDiffusion then uses the trained network to iteratively refine random initial sequences until they match the target structure. We present a comprehensive evaluation and analysis of the proposed method.

Han Huang et al., **RiboDiffusion**: tertiary structure-based RNA inverse folding with generative diffusion models, Bioinformatics, (2024)

I believe that incorporating the backfolded sequence allows the language model to learn key base-pairing dependencies underlying RNA tertiary structure.

# OUTLINE

## ❖ Objectives

## ❖ Dataset

- What is RNA molecule
- RNACentral database

## ❖ Tools description

- RNA alignment Tool
- RNA 3D structure parameters Tool
- The structure-guided feedback Tool

## ❖ LLM Tool-Agent Computational Pipeline

## ❖ Results for LLM Tool-Agent

- Speculative simulation of reverse evolution of tRNA<sup>Gly</sup> (6 slides)

## ❖ LLM ReAct-Agent Computational Pipeline

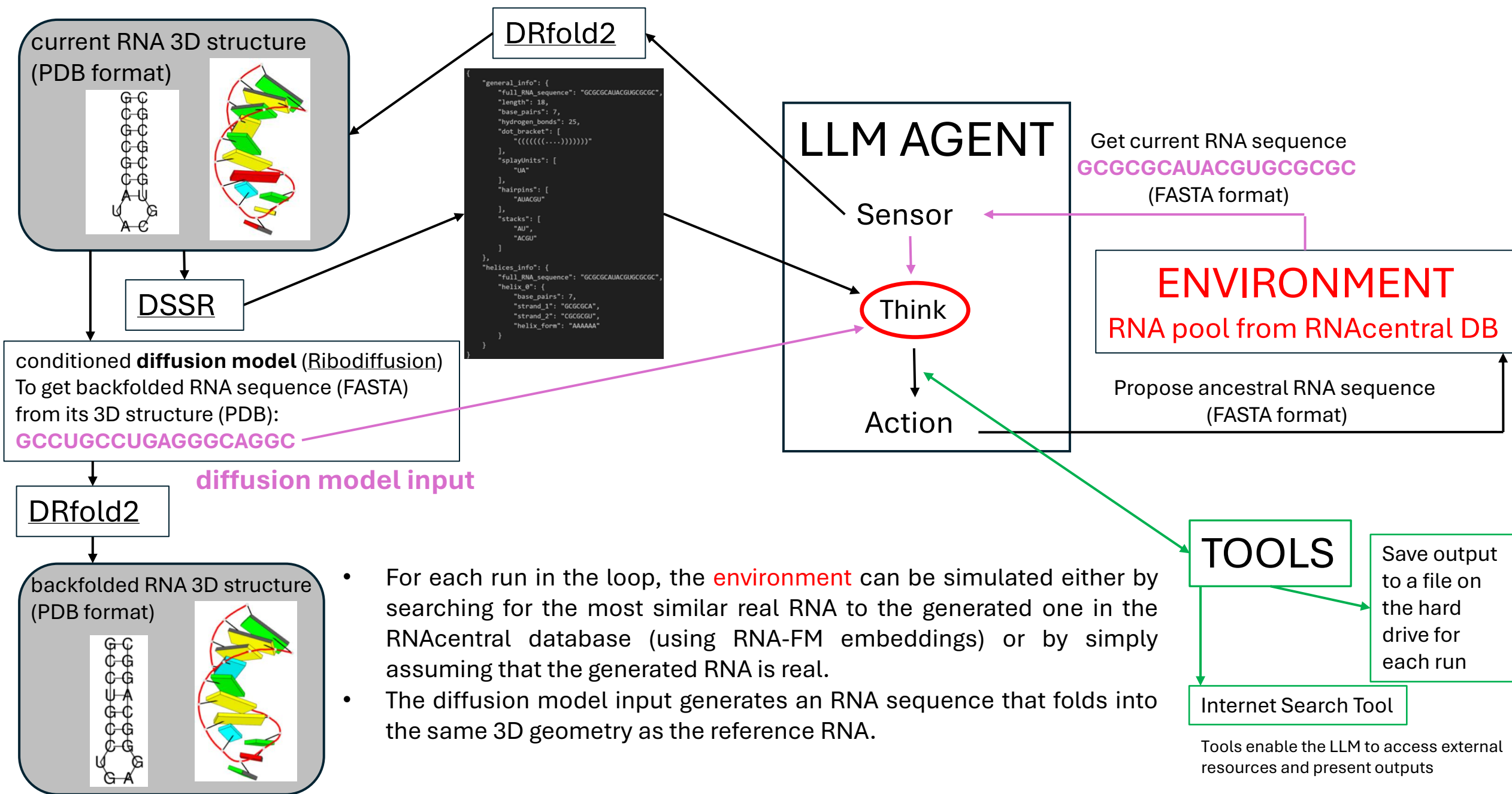
## ❖ Results for LLM ReAct-Agent

- Comparison of two RNAs
- Single point mutation in tRNA associated with MELAS syndrome
- Generate two connected RNA hairpins without RNA template
- Generate two connected RNA hairpins from a given RNA template
- Targeted modification of a specific region of an RNA molecule

## ❖ Conclusions

## ❖ What's next ?

# LLM Tool-Agent Computational Pipeline





# Prompt Engineering: Designing Effective Inputs for LLMs

```
prompt = ChatPromptTemplate.from_messages(  
    [  
        (  
            "system",  
            """  
            You are a research assistant that will help generate a research paper.  
            Answer the user query and use necessary tools.  
            Wrap the output in this format and provide no other text:  
            {format_instructions}  
            Do not include any other commentary or explanation — output only the JSON.  
            """,  
        ),  
        ("placeholder", "{chat_history}"),  
        (  
            "human",  
            """{query}""",  
        ),  
    ],  
    partial(format_instructions=parser.get_format_instructions())
```

Current RNA sequence:  
{current\_sequence}

3D Structural Parameters of the current RNA:  
{rna\_structure\_block}

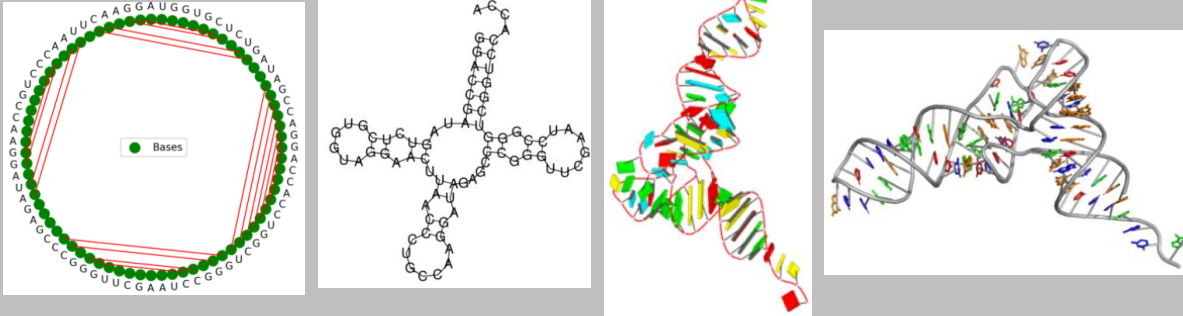
Backfolded RNA sequence from 3D structure of the current RNA:  
{backfolded\_sequence}

```
query = (  
    f"Propose a plausible ancestral RNA sequence based on the current RNA sequence ant its properties."  
    f"Save both a summary and the proposed ancestral_sequence to a file"  
    f"and begin each output in the file with the following section header: evolution_step_{i}"  
    raw_response = agent_executor.invoke({"query": query,  
        "current_sequence": current_sequence,  
        "backfolded_sequence": backfolded_sequence,  
        "rna_structure_block": rna_structure_block})
```

No conditions were imposed on the generated ancestral RNA sequence (e.g., limits on length or number of mutation events per step)!

## STEP 1

GGACCGAUAGUCUCGUGGUAGGAACUUAACCCUGCCAAGG  
AUAGAGCCCGGGUUCGAAUCCGGGUCGGUCCACCA



## DSSR output

```

length: 75
base_pairs: 29
hydrogen_bonds: 112
dot_bracket: ((((((((((.....)))))).((.....))))).(((.....)))))).....
multiplets: UUG, CGG
splayUnits: AU, UG, GUA, CC, AA, CA
ssSegments: ACCA
junctions: AUAGCUUAGAGCGGU
hairpins: CGUGGUAG, CCUGCCAAG, GUUCGAUUC
nonStack: GUUUGA
stacks: AC, GG, GU, GC, GC, AC, GU, UCU, UGG, GAA, AAC, CCU, GAU, GGG, GUU, GACC, AAG, ACAU, CAAg
non_coxStacks: 2
helix 0: base_pairs=12, strand_1=GGACGACCGGG, strand_2=CUUGGUUGGGCC, helix_form=AAAAAXAA.A
helix 1: base_pairs=13, strand_1=CAGGAUAGUCUG, strand_2=UCCCAAUCAAGGA, helix_form=X.A.AAA.X

```

The provided 75-nucleotide RNA sequence exhibits characteristic tRNA-like structural features with a cloverleaf secondary structure. Key structural parameters include:

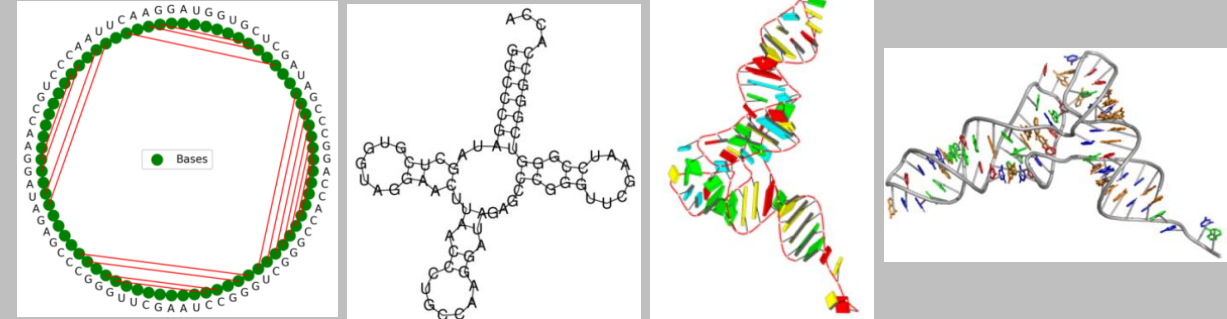
- Length: 75 nucleotides
- Base pairs: 29
- Dot-bracket notation showing complex stem-loop structures
- Multiple hairpins and junction regions
- Two coaxial stacks indicating sophisticated tertiary interactions

The first two images in each panel were generated using simplified 2D structures for rough visualization purposes only.

The last two images depict the actual 3D structures, predicted by DRfold2 and further processed by DSSR to extract structural parameters for LLM input.

## ANCESTRAL\_SEQUENCE:

GGCCCGAUAGCUCGUGGUAGGAACUUAACCCUGCCAAGGAUAGAG  
CCCGGGUUCGAAUCCGGGUCGGGCCACCA



## ## Rationale for Ancestral Sequence Reconstruction

1. **\*\*Structural Conservation\*\***: The ancestral sequence maintains the essential cloverleaf structure characteristic of early tRNA-like molecules, preserving critical base-pairing regions that ensure functional stability.
2. **\*\*Reduced Complexity\*\***: Early RNA molecules likely had simpler tertiary structures with fewer modified nucleotides and less sophisticated coaxial stacking arrangements.
3. **\*\*Enhanced Thermostability\*\***: Ancestral RNA sequences typically exhibited greater thermodynamic stability to survive in harsh primordial conditions, achieved through stronger base-pairing interactions.
4. **\*\*Functional Core Preservation\*\***: The proposed ancestral sequence maintains the essential functional domains including the acceptor stem, anticodon loop region, and critical junction sequences necessary for aminoacylation and ribosome binding.
5. **\*\*Evolutionary Parsimony\*\***: The reconstruction follows the principle of minimal evolutionary changes while preserving essential structural and functional characteristics required for early protein synthesis machinery.

The proposed ancestral sequence represents a plausible evolutionary precursor that could have given rise to the current RNA through gradual structural refinements and optimization for specific cellular environments.

Based on analysis of the current RNA sequence (75 nucleotides) with its complex secondary structure containing 29 base pairs, multiple hairpins, and helical regions, I propose an ancestral sequence that would have been simpler yet structurally viable. The current sequence shows sophisticated features including two major helices, multiple hairpin loops, and complex junction regions.

For the ancestral reconstruction, I considered several evolutionary principles:

1. Ancient RNA sequences likely had simpler secondary structures with fewer complex loops
2. Early RNA would have had stronger base-pairing stability in core regions
3. Primitive sequences would have contained more canonical Watson-Crick pairs
4. The ancestral form would retain essential structural motifs but with reduced complexity

The proposed ancestral sequence maintains the basic cloverleaf-like structure but with:

- Simplified hairpin loops with more stable base-pairing
- Reduced junction complexity
- Stronger stem regions with canonical base pairs
- Shorter variable loops
- Enhanced stability through increased GC content in critical structural regions

This ancestral form would represent an evolutionary intermediate that possessed the fundamental structural framework while being simpler and more thermodynamically stable than the modern sequence.

## Current-Ancestral RNA sequences alignment

```

<class 'Bio.Align.Alignment'>
target      0  GGACCGAUAAGUCUGUGGUAAGGAACUUAACCCUGCGCAAGGAUAGAGCCCGGGUUCGAU
               0  | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
query       0  GGCCCGUAUAG-CUCUGUGGUAAGGAACUUAACCCUGCGCAAGGAUAGAGCCCGGGUUCGAU
               0  | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

target      60  CGGGUCGGUCACCA 75
               60  | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
query       59  CGGGUCGGGCGACCA 74

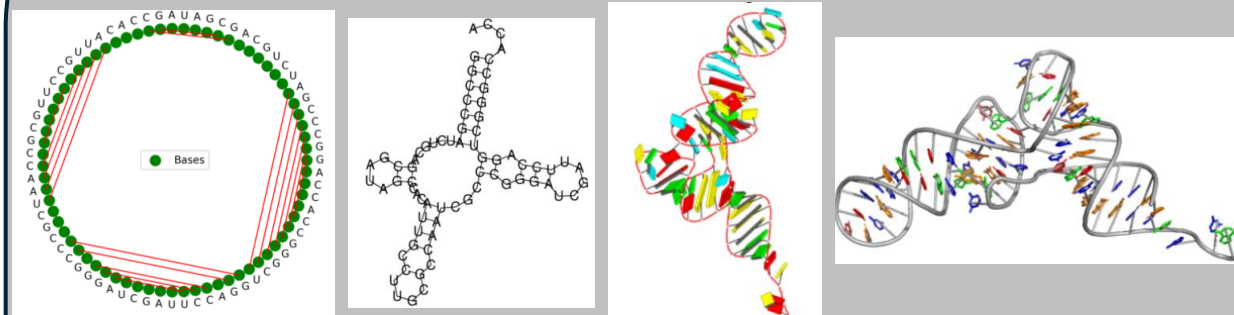
Score: 140.0

```

To improve the reliability of these evolutionary trajectories, it is necessary to introduce constraints on the number of mutation events per step.

## STEP 6

GGCCCGAUCUGCAGCGAUAGCCACAUUGCCUUGCGCCAAUCGCC  
CGGGAUCGAUUCCAGGUCGGGCCACCA



DSSR output

```
length: 71
base_pairs: 32
hydrogen_bonds: 111
dot_bracket: ((((((.....((.....))....(((.....))....))....(((.....))....))....))....
multiplets: CGC, UGCG
splayUnits: UC, GAUA, UU, CC, UU, CC
ssSegments: ACCA
junctions: AUCUGCGCCACAUCGCGU
hairpins: CGAUAG, CUUGCG, AUCGAUU
nonStack: AUUU
stacks: CU, UC, CG, AU, CA, GU, CA, GAC, AUC, GGC, CAG, CGU, GGG, CCAC, GCCU, UGCGCC
num_coaxStacks: 1
helix_0: base_pairs=13, strand_1=GGCCCGACCGGA, strand_2=CGGGGUGGACCU, helix_form=AAAAAxA..AA
helix_1: base_pairs=8, strand_1=UGAUGCU, strand_2=GUACGC, helix_form=xAAAx..
helix_2: base_pairs=5, strand_1=CAGCG, strand_2=ACCGC, helix_form=..Ax
```

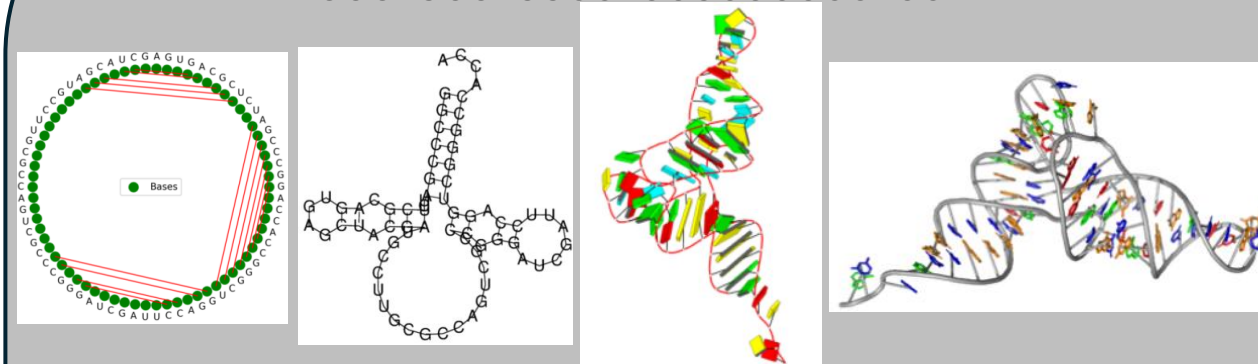
Based on the analysis of the current RNA sequence (71 nucleotides) and its structural properties, I propose an ancestral RNA sequence that maintains the critical structural elements while reflecting probable evolutionary simplifications. The current RNA displays a complex secondary structure with 32 base pairs, 3 major helices, and multiple structural motifs including hairpins, bulges, and junctions.

### Key observations from the current sequence:

- Length: 71 nucleotides
- Three major helices (13, 8, and 5 base pairs respectively)
- Complex dot-bracket notation indicating nested stem-loop structures
- Multiple hairpin loops (CGAUAG, CUUGCG, AUCGAUU)
- Junction regions and single-stranded segments
- Strong GC content in stem regions providing stability

## ANCESTRAL\_SEQUENCE:

GGCCCGAUCUCGCAGUGAGCUACGAUGCCUUGCGCCAGUCGCCC  
GGGAUCGAUUCCAGGUCGGGCCACCA



The proposed ancestral sequence likely represents a simpler evolutionary precursor that:

1. Maintains the core three-helix architecture
2. Has simplified loop regions with fewer complex motifs
3. Retains essential base-pairing patterns for structural stability
4. Shows reduced sequence complexity while preserving functional domains
5. Exhibits characteristics consistent with early RNA world molecules

### The ancestral reconstruction considers:

- Conservation of critical structural elements
- Simplification of complex motifs that likely evolved later
- Maintenance of thermodynamic stability
- Preservation of potential catalytic or binding sites
- Evolutionary parsimony in sequence changes

```
<class 'Bio.Align.Alignment'>
target      0 GGCCCGAUCU-GCAGCGAUAGCCACAUUGCCUUGCGCCAAUCGCCCGGGAUCGAUCCAG
              0 |||||...-|||...-|||...-|||...-|||...-|||...-|||...-|||...-|||...
query       0 GGCCCGAUCUCGAGUG--AGCUACGAUGCCUUGCGCCAGUCGCCCGGGAUCGAUCCAG

target      59 GUCGGGCCACCA 71
              60 |||||... 72
query       58 GUCGGGCCACCA 70

Score: 118.5
```

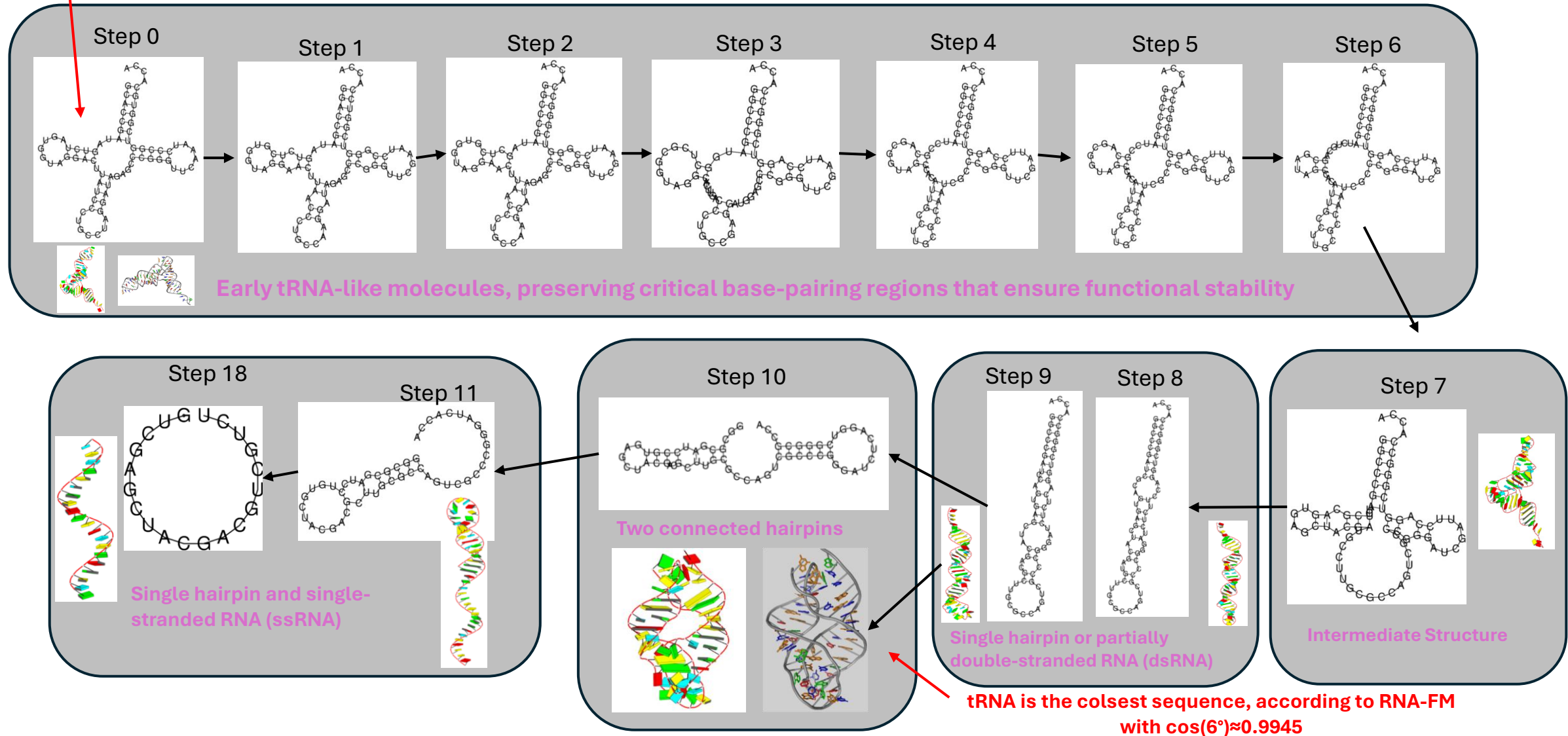
The first two images in each panel were generated using simplified 2D structures for rough visualization purposes only.

The last two images depict the actual 3D structures, predicted by DRfold2 and further processed by DSSR to extract structural parameters for LLM input

To improve the reliability of these evolutionary trajectories, it is necessary to introduce constraints on the number of mutation events per step.

# Results of Simulated Reverse Evolution of tRNA<sup>Gly</sup> Asgard\_group\_archaeon (First Run)

tRNA<sup>Gly</sup>  
Asgard\_group\_archaeon  
Starting point

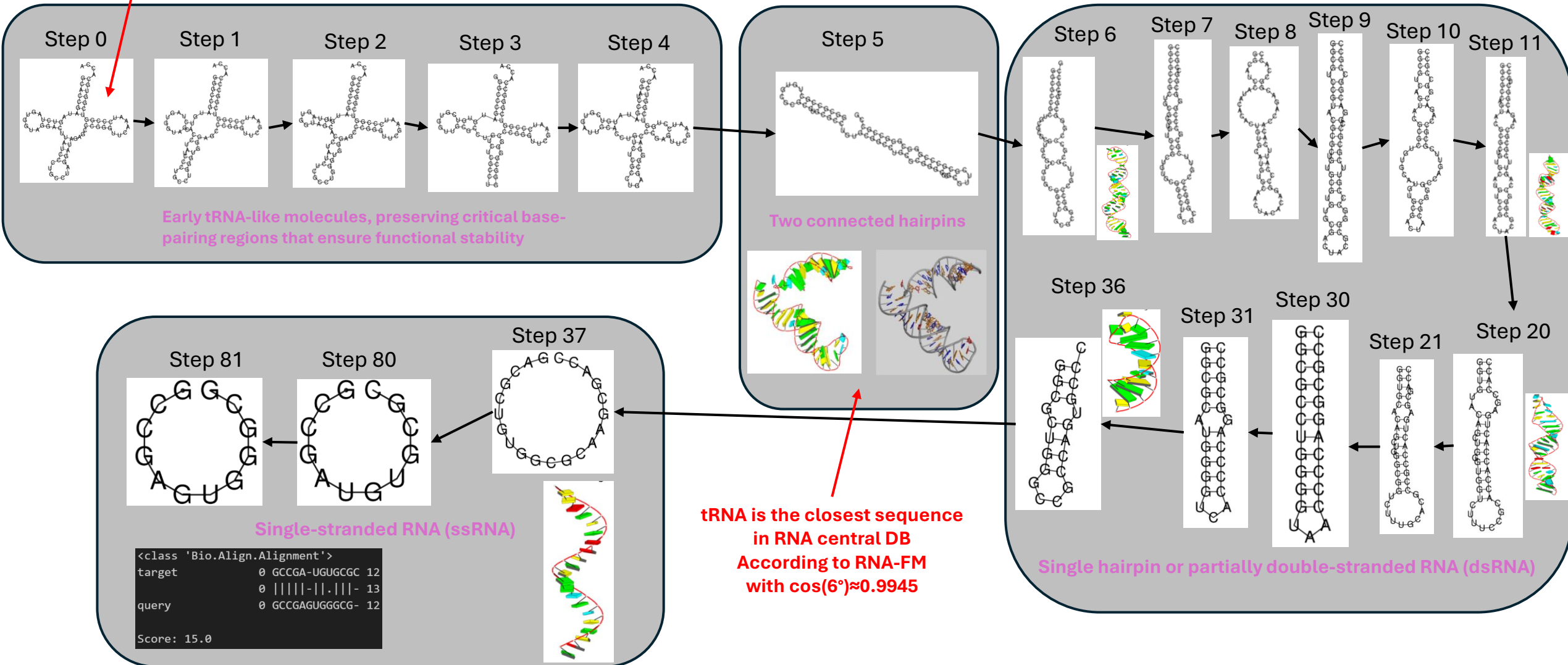


To improve the reliability of these evolutionary trajectories, it is necessary to introduce constraints on the number of mutation events per step.



# Results of Simulated Reverse Evolution of tRNA<sup>Gly</sup> Asgard\_group\_archaeon (Second Run)

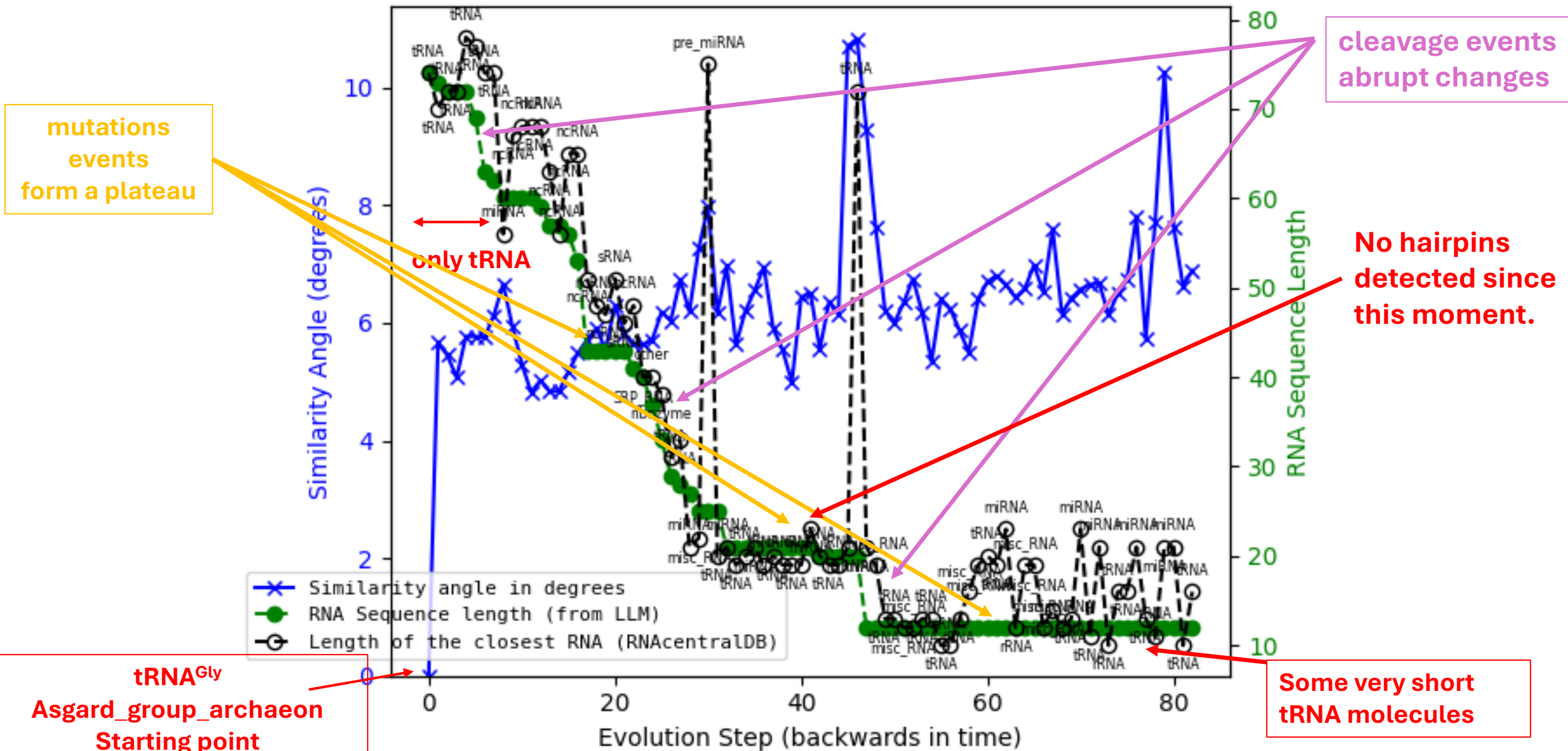
tRNA<sup>Gly</sup>  
Asgard\_group\_archaeon  
Starting point



To improve the reliability of these evolutionary trajectories, it is necessary to introduce constraints on the number of mutation events per step.

# Using Large Language Models to Simulate the Reverse Evolution of tRNA<sup>Gly</sup>

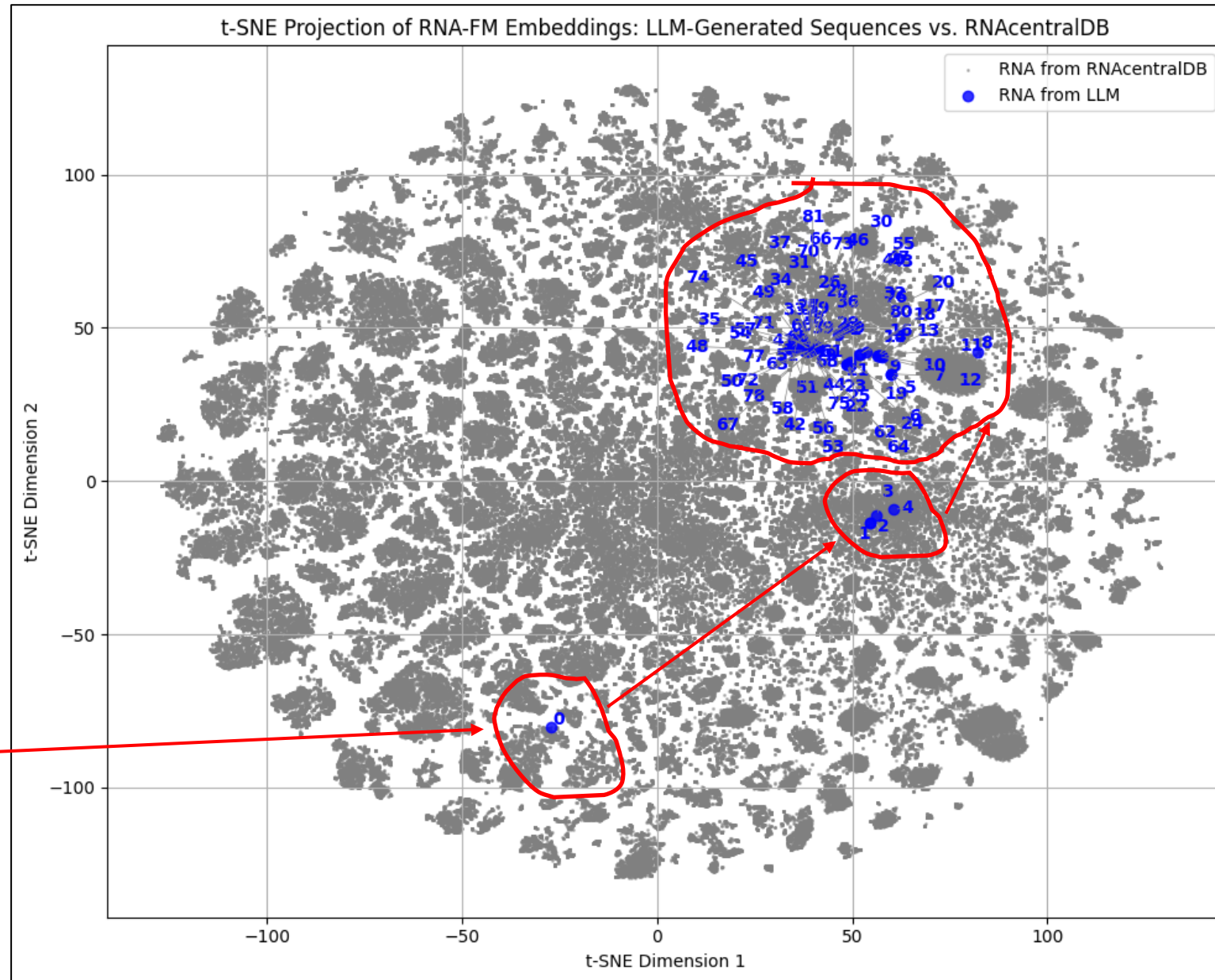
Smallest similarity angle between RNAs from LLM and RNACentralDB (comparison using RNA-FM)



From tRNA<sup>Gly</sup> to short tRNA and miscellaneous RNA



# How Accurately do LLM-Designed RNAs Represent Real RNA?



0, 1, ... are  
evoution step  
(backwards in  
time)

tRNA<sup>Gly</sup>  
Asgard\_group\_archaeon  
Starting point

400,000 tRNA sequences from the RNAcentral database with lengths less than 80 nucleotides

# OUTLINE

## ❖ Objectives

## ❖ Dataset

- What is RNA molecule
- RNACentral database

## ❖ Tools description

- RNA alignment Tool
- RNA 3D structure parameters Tool
- The structure-guided feedback Tool

## ❖ LLM **Tool-Agent** Computational Pipeline

## ❖ Results for LLM Tool-Agent

- Speculative simulation of reverse evolution of tRNA<sup>Gly</sup> (6 slides)

## ❖ LLM **ReAct-Agent** Computational Pipeline

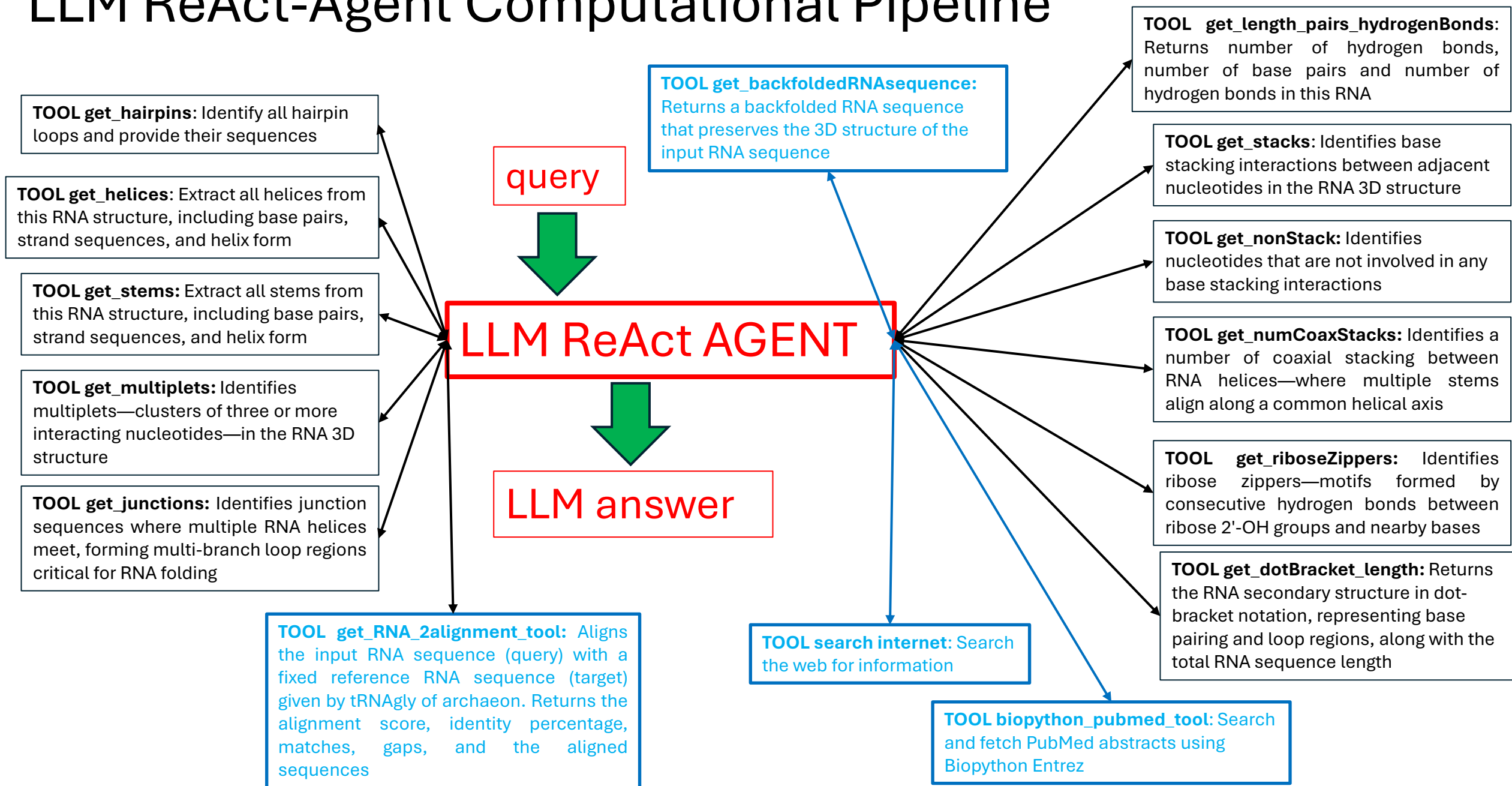
## ❖ Results for LLM ReAct-Agent

- Comparison of two RNAs
- Single point mutation in tRNA associated with MELAS syndrome
- Generate two connected RNA hairpins without RNA template
- Generate two connected RNA hairpins from a given RNA template
- Targeted modification of a specific region of an RNA molecule

## ❖ Conclusions

## ❖ What's next ?

# LLM ReAct-Agent Computational Pipeline



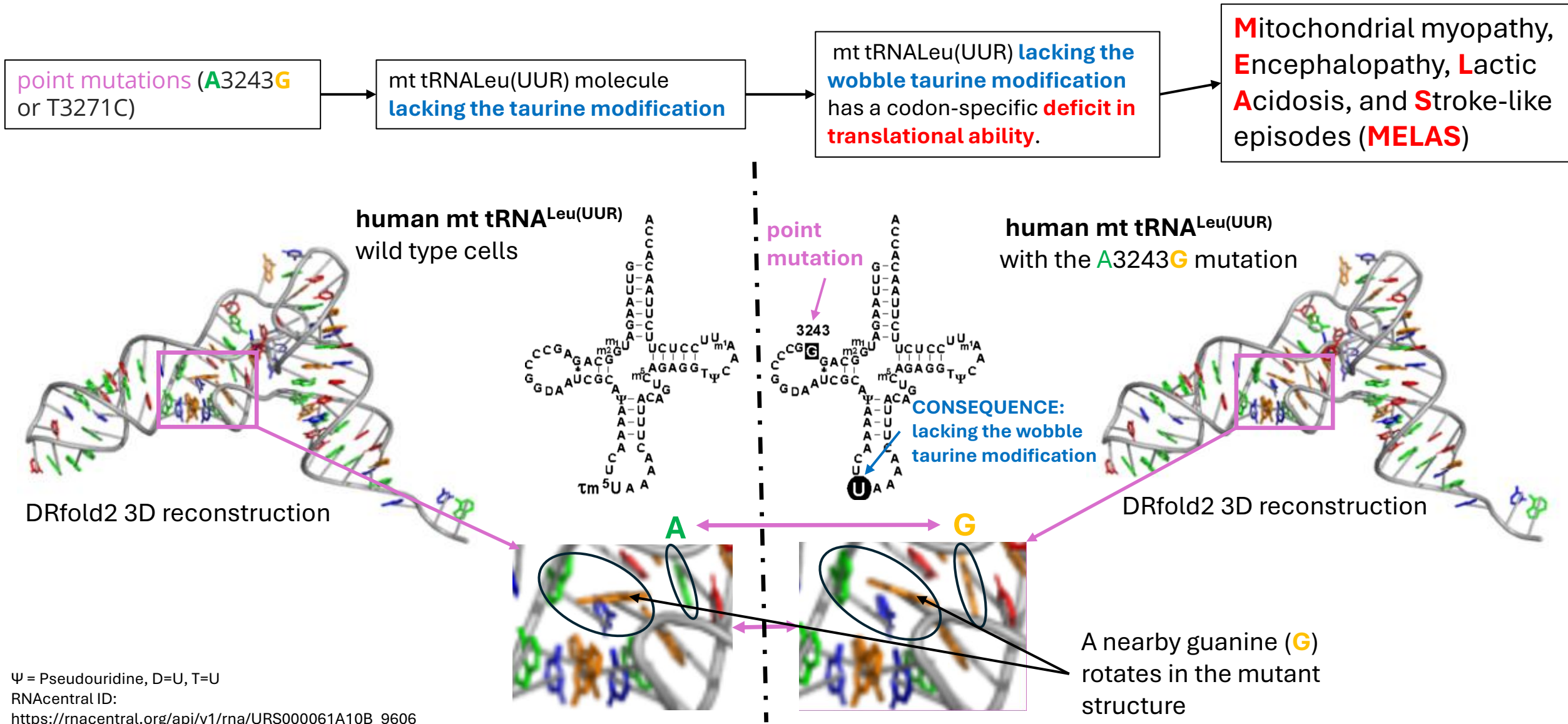
## Agent: Comparison of Two RNA Sequences



Provide the output in clearly labeled sections.

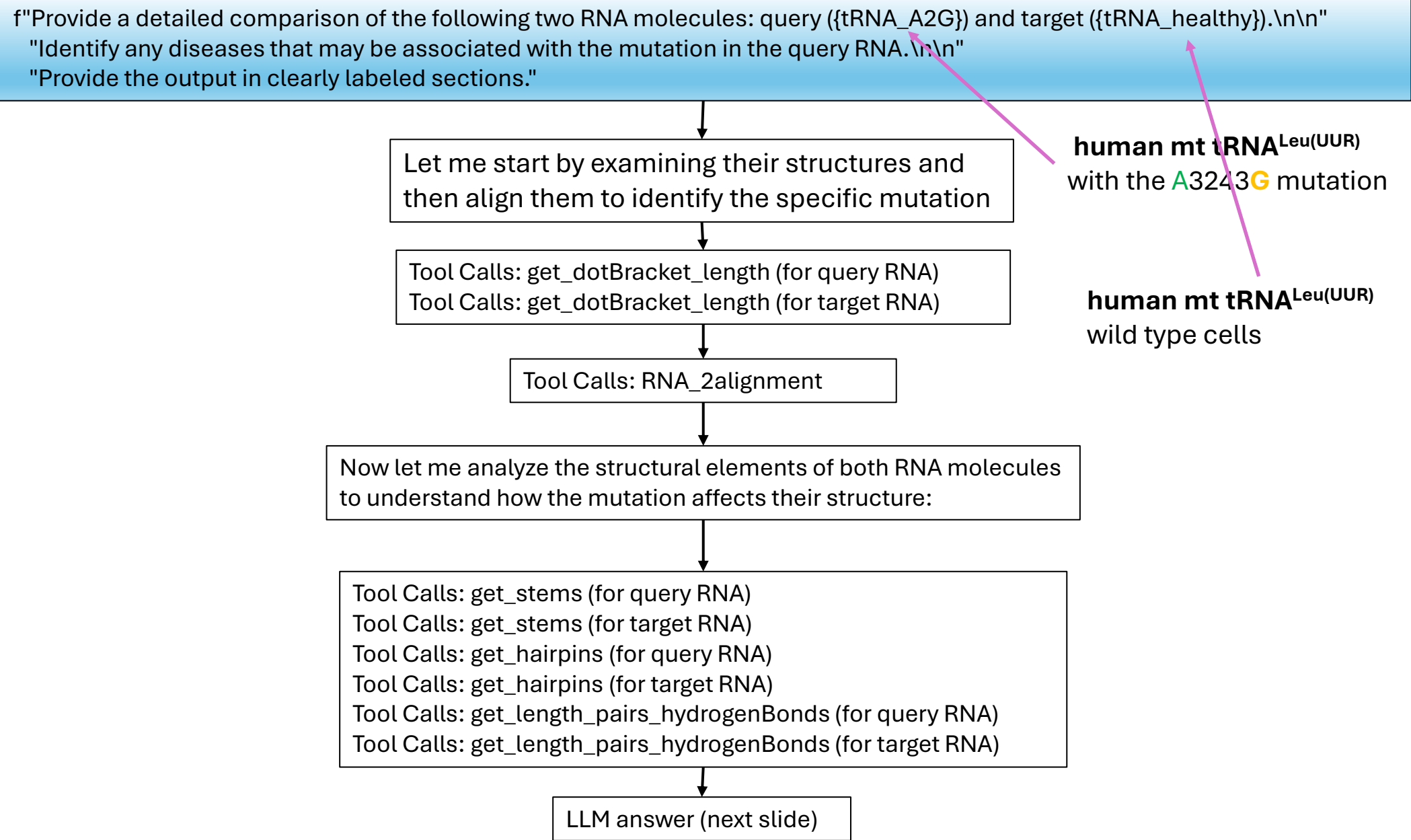
Aspect	Query (AG repeat)	Target (tRNA-like)
**Length**	20 nt	74 nt
**Structure**	Completely unstructured	Highly structured (33 base pairs)
**Complexity**	Simple alternating repeat	Complex with 4 helices, 3 hairpins
**Function**	Synthetic/artificial	Functional tRNA-like molecule
**Stability**	Low (no base pairs)	High (123 hydrogen bonds)
**Evolution**	Likely artificial	Biologically evolved structure
**Diversity**		

# Single Point Mutation in tRNA Associated with MELAS Syndrome





# Single Point Mutation in tRNA Associated with MELAS Syndrome





# Single Point Mutation in tRNA Associated with MELAS Syndrome

# \*\*Detailed RNA Comparison Analysis\*\*

## ## \*\*1. Sequence Alignment Results\*\*

### ### \*\*Basic Statistics:\*\*

- \*\*Sequence Length:\*\* 78 nucleotides (both sequences)
- \*\*Alignment Score:\*\* 153.0
- \*\*Identity:\*\* 98.72% (77/78 matches)
- \*\*Mismatches:\*\* 1
- \*\*Gaps:\*\* 0 in both sequences

### ### \*\*Mutation Identification:\*\*

- \*\*Position:\*\* 14 (counting from 5' end)
- \*\*Mutation Type:\*\* G→A transition
- \*\*Query sequence:\*\*

GUUAAGAUGGCAG\*\*G\*\*GCCCCGUAAUCGCAUAAAACUAAAACUUUACAGUCAGAGGUUCAUUCCUC  
UUCUUAACACCA

- \*\*Target sequence:\*\*

GUUAAGAUGGCAG\*\*A\*\*GCCCCGUAAUCGCAUAAAACUAAAACUUUACAGUCAGAGGUUCAUUCCUC  
UUCUUAACACCA

## ## \*\*2. Secondary Structure Comparison\*\*

### ### \*\*Dot-Bracket Notation:\*\*

Both RNA molecules have identical secondary structure patterns:

```\n

(((((.....((.....)).(((.....))))).(((.....)))))).....(((.....)))))).....\n```\n

### ### \*\*Structural Elements:\*\*

\*\*Query RNA (with G at position 14):\*\*

- \*\*Stems:\*\* 4 helical regions with identical base pairing patterns
- \*\*Hairpin Loops:\*\*
  - Loop 1: CAGGCCCCGUAAUCG (16 nucleotides)
  - Loop 2: AACUAAAACU (11 nucleotides)
  - Loop 3: GUUCAAUUC (9 nucleotides)
- \*\*Base Pairs:\*\* 34
- \*\*Hydrogen Bonds:\*\* 106

\*\*Target RNA (with A at position 14):\*\*

- \*\*Stems:\*\* 4 helical regions with identical base pairing patterns
- \*\*Hairpin Loops:\*\*
  - Loop 1: CAGAGCCCCGUAAUCG (16 nucleotides)
  - Loop 2: AACUAAAACU (11 nucleotides)
  - Loop 3: GUUCAAUUC (9 nucleotides)
- \*\*Base Pairs:\*\* 37
- \*\*Hydrogen Bonds:\*\* 108

## ## \*\*3. Structural Impact Analysis\*\*

### ### \*\*Key Differences:\*\*

1. \*\*Base Pairs:\*\* Target RNA has 3 additional base pairs (37 vs 34)
2. \*\*Hydrogen Bonds:\*\* Target RNA has 2 additional hydrogen bonds (108 vs 106)
3. \*\*Loop Sequence:\*\* The first hairpin loop differs by one nucleotide (GGG vs GAG)

### ### \*\*Structural Stability:\*\*

The G→A mutation appears to **increase structural stability** in the target RNA by:

- Creating additional base pairing opportunities
- Increasing the total number of hydrogen bonds
- Potentially stabilizing the first hairpin loop region

## ## \*\*4. RNA Identity and Disease Associations\*\*

**Based on the sequence characteristics and structure, this RNA appears to be a **transfer RNA (tRNA)** molecule, given its:**

- **Length (~78 nucleotides)**
- **Secondary structure with multiple stems and loops**
- **Characteristic tRNA-like folding pattern**

### ### \*\*Potential Disease Associations:\*\*

**G→A transitions in tRNA molecules** are commonly associated with:

## 1. **Mitochondrial Diseases:**

- **MELAS (Mitochondrial Encephalomyopathy, Lactic Acidosis, and Stroke-like episodes)**

- MERRF (Myoclonus Epilepsy with Ragged Red Fibers)
- Various mitochondrial myopathies

## 2. **Metabolic Disorders:**

- Respiratory chain deficiencies
- Energy metabolism dysfunction

## 3. **Neurological Conditions:**

- Progressive external ophthalmoplegia
- Leigh syndrome variants

**Note:** The specific disease association would depend on:

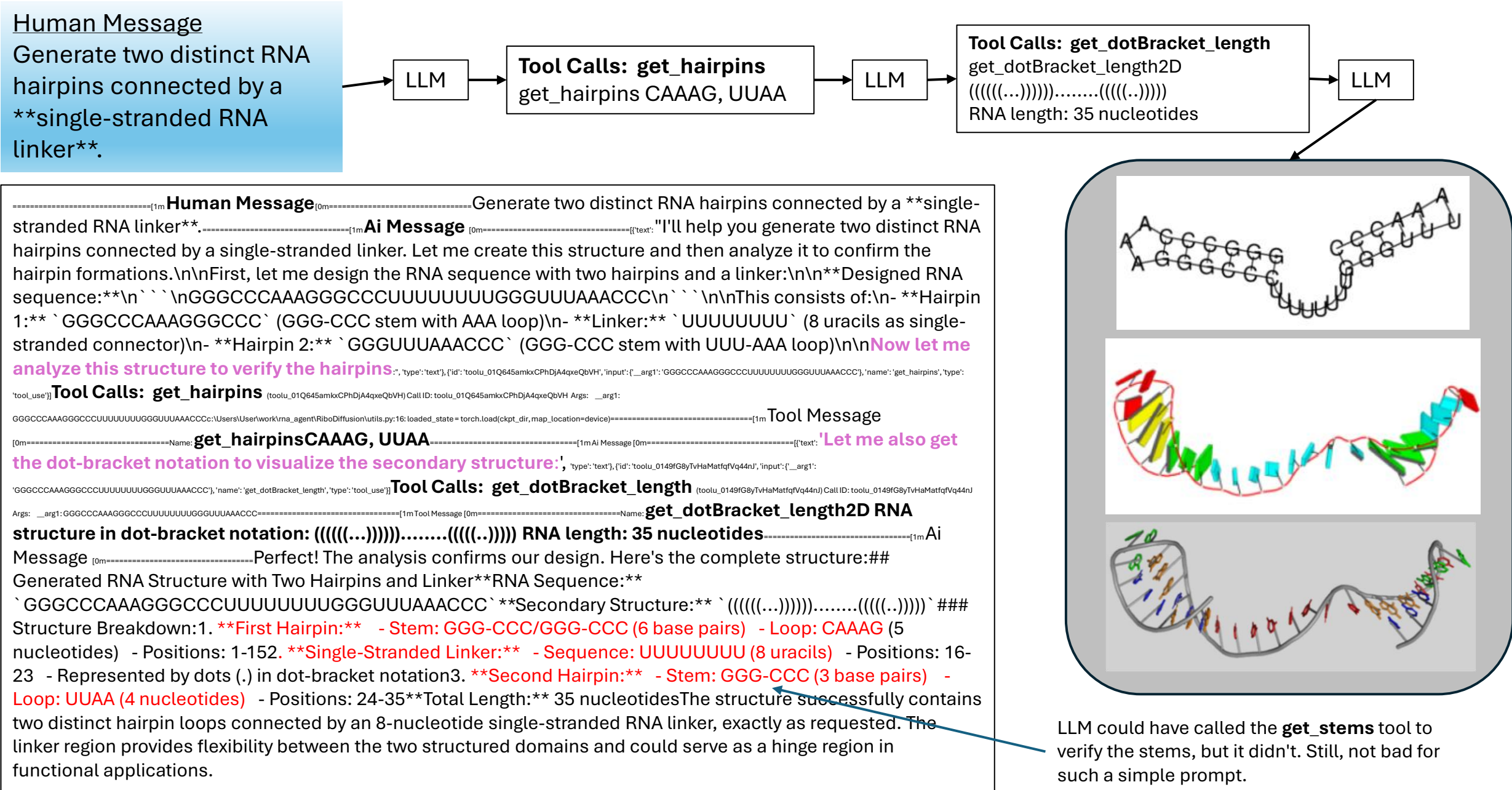
- Which specific tRNA this represents
- The cellular location (mitochondrial vs cytoplasmic)
- The functional impact on amino acid charging or translation

## ## \*\*5. Functional Implications\*\*

The mutation from G to A at position 14:

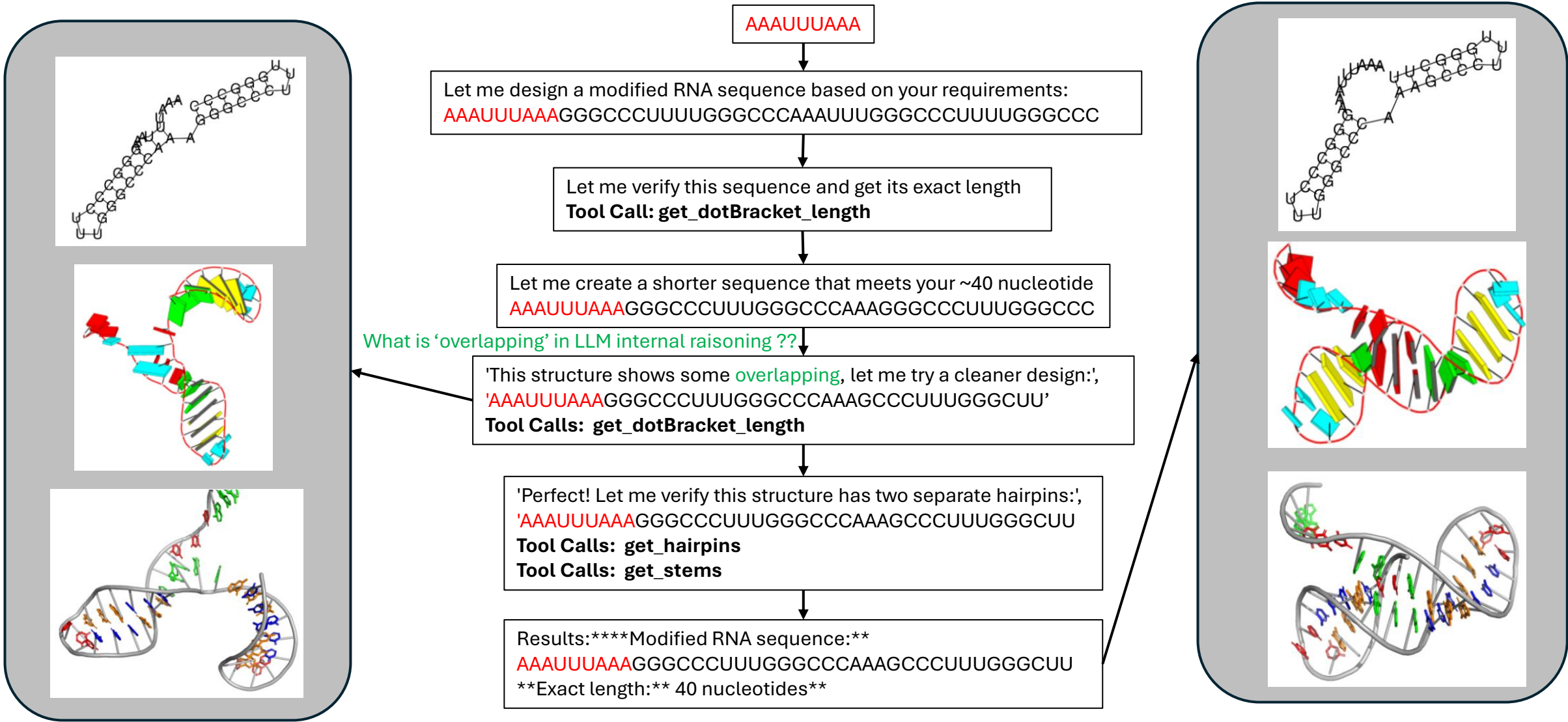
- **Increases structural stability** (more base pairs and H-bonds)

# LLM ReAct agent : Generate Two Connected RNA Hairpins without a Prior RNA Template



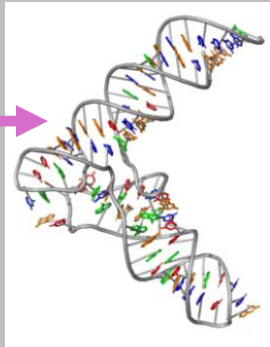
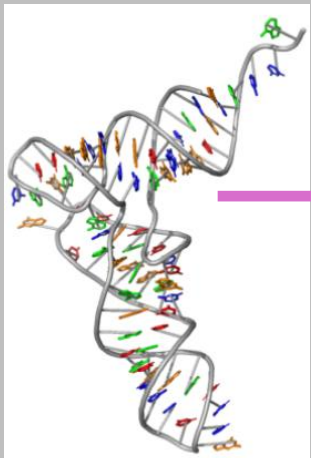
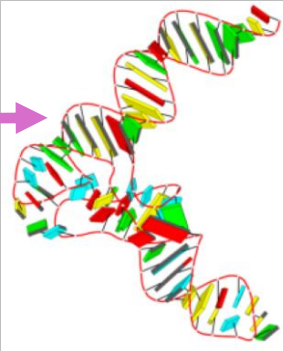
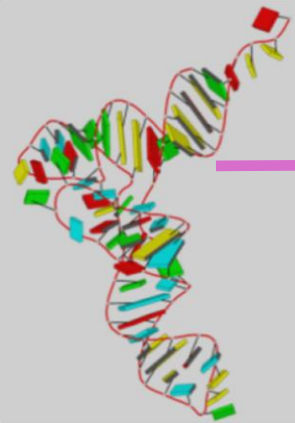
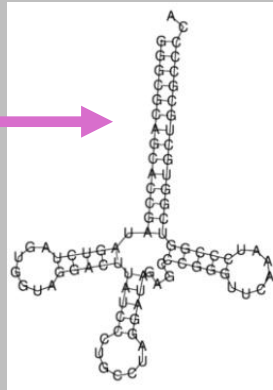
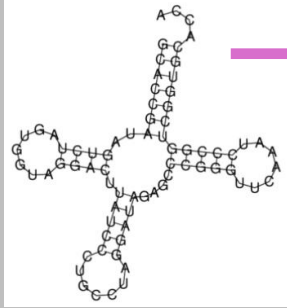
# LLM ReAct agent : Generate Two Connected RNA Hairpins with a Prior RNA Template

The reference sequence (input): **AAAUUUAAA**  
Modify this reference sequence to form a structure without **\*\*pseudoknot base pairs\*\*** consisting of **\*\*two separate hairpins\*\*** connected by a **\*\*single-stranded RNA linker\*\***.  
The modified RNA sequence must be about 40 nucleotides long. Return the modified RNA sequence and its exact length.



# reference sequence      modified sequence

tRNA<sup>Gly</sup>  
Asgard\_group\_archaeon



## LLM ReAct Agent: Targeted modification of a specific region of an RNA 1

f"Current tRNA sequence (input): {tRNA}\n\n"

"Modify this reference sequence by increasing the length of the acceptor stem to **\*\*twice its original size\*\***. "

"Ensure that **\*\*no more than 20%\*\*** of the nucleotides are changed or added.\n\n"

"Keep the rest of the tRNA unchanged.\n\n"

"Verify the structure and return the modified RNA sequence."

Current tRNA sequence (input):

GCACCGAUAGUCUAGUGGUAGGACUUAUCCUGCCUAGGAUAGAGCCCGGGUCAAUCCCGGUCGGUGCACCA

I'll help you modify the tRNA sequence by doubling the acceptor stem length while keeping changes to no more than 20% of nucleotides. First, let me analyze the current structure to identify the acceptor

Tool Calls: **get\_dotBracket\_length**

'Now let me identify the stems to locate the acceptor stem

Tool Calls: **get\_stems**

'Let me also get the hairpin loops to better understand the structure

Tool Calls: **get\_hairpins**

Based on the analysis, I can see this is a standard tRNA structure with:\n- **\*\*Acceptor stem\*\***: **helix\_0 with 7 base pairs (GCACCGA paired with CGUGGCU)**\n- **D-arm**: **helix\_1 with 4 base pairs**\n- **Anticodon arm**: **helix\_2 with 5 base pairs** \n- **TψC arm**: **helix\_3 with 5 base pairs**\n\nTo double the acceptor stem length from 7 to 14 base pairs while changing ≤20% of nucleotides (≤15 out of 74), I'll:\n1. Extend the 5' end with 7 complementary base pairs\n2. Extend the 3' end with 7 complementary base pairs\n3. Keep all other structural elements unchanged\n\nHere's the modified sequence: **GGGC...**

Tool Calls: **get\_dotBracket\_length**

Tool Calls: **get\_stems**

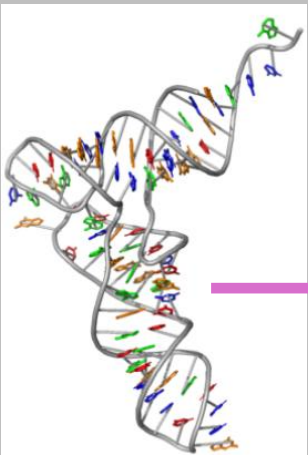
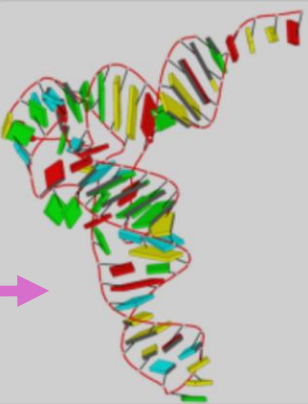
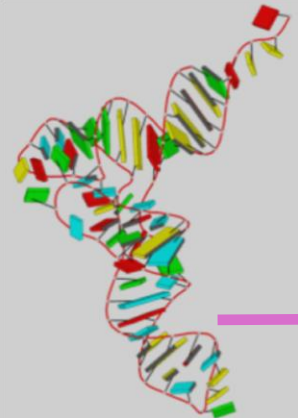
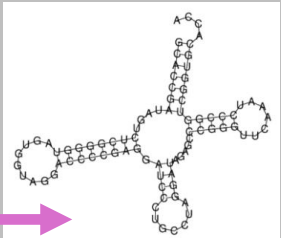
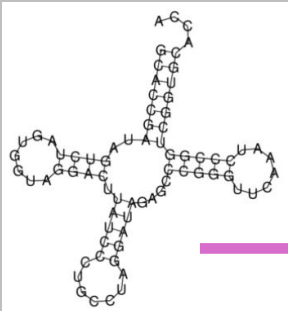
Figure adapted from Frontiers in Molecular Biosciences (Irem Avcilar-Kucukgoze et al., 2020), CC BY 4.0. DOI: 10.3389/fmolb.2020.610617

Perfect! The modification is successful:## Summary:- **\*\*Original length\*\***: 74 nucleotides- **\*\*Modified length\*\***: 86 nucleotides (12 nucleotides added = 16.2% increase, well under 20% limit)- **\*\*Acceptor stem\*\***: **Successfully doubled from 7 to 14 base pairs**- **\*\*Structure integrity\*\***: All other arms (D-arm, anticodon arm, TψC arm) remain unchanged- **\*\*Base pair complementarity\*\***: Maintained throughout the extended acceptor stem## Modified tRNA sequence:\*\*`GGGCGCAGCACCGAUAGUCUAGUGGUAGGACUUAUCCUGCCUAGGAUAGAGCCCGGGUCAAUCCCGGUCGGUGCUGCGCCCCA`\*\*The acceptor stem is now twice its original length (14 bp vs 7 bp) while preserving the overall tRNA secondary structure and staying within the 20% nucleotide change constraint.



reference sequence modified sequence

tRNA<sup>Gly</sup>  
Asgard\_group\_archaeon



LLM ReAct Agent: Targeted modification of a specific region of an RNA 2

f"Current tRNA sequence (input): {tRNA}\n\n"

"Modify this reference sequence by increasing the length of the D-stem to **\*\*twice its original size\*\***. "

"Use **backfolded sequence** which forms the same 3D structure. \n\n"

"Keep the rest of the tRNA sequence and its cloverleaf structure unchanged.\n\n"

"Verify the structure and return the modified RNA sequence."

It does not  
work  
without  
backfolded  
sequence

Let me start by analyzing the current structure to identify the D-stem and then create the modified sequence.

Tool Calls: get\_stems  
get\_dotBracket\_length  
get\_backfoldedRNAsequence

Looking at the structure, the D-stem is located at positions approximately 10-13 and 22-25 in the original sequence. To double the D-stem length while preserving the overall tRNA structure, I'll extend both strands of the D-stem by 4 additional base pairs. \n\nHere's my strategy: \n1. Keep the acceptor stem, anticodon stem, and TψC stem unchanged \n2. Extend the D-stem from 4 to 8 base pairs by adding complementary sequences \n3. Use the backfolded approach to maintain 3D structure \n\nLet me create the modified sequence by extending the D-stem: \n\nOriginal D-stem: GUCU...CAGG (4 bp) \nExtended D-stem: GUCUGGCC...CCAGGACG (8 bp) \n\nThe modified tRNA sequence would be: \nGCACCGAUAGUCUGGCCAGUGGUAGGACGCCAGGAUCCUGCCUAGGAUAGAGCCCGGGUCAAUCCCGGUCGGUGCACCA

Tool Calls: get\_stems

get\_dotBracket\_length

'Good! The D-stem is now 5 base pairs, but I need it to be 8 (double the original 4). Let me create a better modification: \n\nThe modified sequence with D-stem extended to 8 base pairs: \nGCACCGAUAGUCUCGGGCAGUGGUAGGACGCCGAGGAUCCUGCCUAGGAUAGAGCCCGGGUCAAUCCCGGUCGGUGCACCA'

Tool Calls: get\_stems

'Getting closer! Let me make one more adjustment to achieve exactly 8 base pairs in the D-stem: \n\nGCACCGAUAGUCUCGGGGAGUGGUAGGACCCCCGAGGAUCCUGCCUAGGAUAGAGCCCGGGUCAAUCCCGGUCGGUGCACCA'

Tool Calls: get\_stems

Let me try a more precise approach to get exactly 8 base pairs: \n\nGCACCGAUAGUCUCGGGGUAGUGGUAGGACCCCCGAGGAUCCUGCCUAGGAUAGAGCCCGGGUCAAUCCCGGUCGGUGCACCA'

Tool Calls: get\_stems

'Perfect! Now I have the D-stem with exactly 8 base pairs (double the original 4). Let me get the backfolded sequence and verify the complete structure: \n\nTool Calls: get\_backfoldedRNAsequence  
GCCCCAGUAGAGGCCGGCUUAGGUGUCACCCGACCGGUGGUCAGACCCAUGACCCACGGGUUAAAGUCCCGUCUGGGGCACCA

GCACCGAUAGUCUCGGGGUAGUGGUAGGACCCCCGAGGAUCCUGCCUAGGAUAGAGCCCGGGUCAAUCCCGGUCGGUGCACCA  
===== Tool Message =====  
Name: get\_dotBracket\_length

I have successfully modified the tRNA sequence by doubling the D-stem length from 4 to 8 base pairs.

# CONCLUSIONS

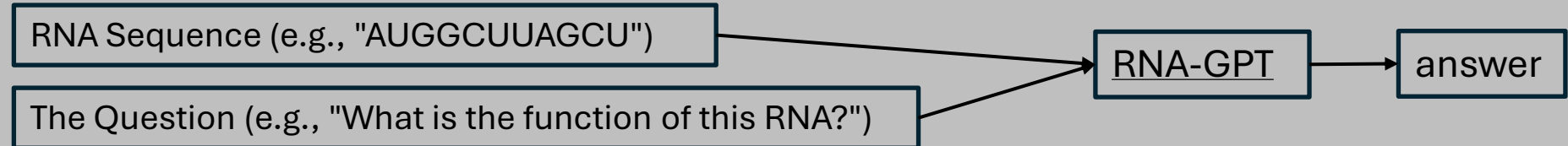
- ❖ **An LLM agent can be used not only to generate new RNA sequences but also to modify existing ones, offering an advantage over conditional diffusion models alone when guiding RNA design.**
- ❖ **Without external tools, an LLM is not even able to reliably count the number of nucleotides in an RNA sequence. A backfolded sequence turns out to be an effective, LLM-comprehensible alternative to 3D structural parameters for preserving the original structure of an input (template) sequence during partial modification.**
- ❖ **A preliminary, speculative attempt to model tRNA reverse evolution captures certain trends. RNA length decreases over time, and three distinct evolutionary patterns have been identified:**
  - *The ancestral sequence maintains the essential cloverleaf structure characteristic of early tRNA-like molecules, preserving critical base-pairing regions.*
  - *A cloverleaf structure transforms into a single hairpin or partially double-stranded RNA (dsRNA) through intermediate structures (e.g., two connected hairpins).*
  - *Finally, the hairpin or dsRNA transforms into single-stranded RNA (ssRNA).*



# 1. What's next ?

1. Introduce **limits** on the length or **the number of mutation events per step**
2. Search for mistakes in the answers and correct the prompt accordingly
3. Comprehensive literature research
4. Add more Tools

## 5. Use RNA-GPT to get more context from answering questions about given RNA sequences



The RNA-FM sequence encoder, which accounts for 3D structure, was used to embed RNA sequences for alignment with natural language, at this time, RNA-GPT is not available for use

Yijia Xiao et al. , RNA-GPT: Multimodal Generative System for RNA Sequence Understanding, **2025**

## 6. Fundamental Studies ?

- The oldest part of ribosome: pseudosymmetrical region (~90+90nt)
- The Spiegelman monster: the shortest “self-replicating” RNA molecule (~220nt)

large ribosomal subunit (LSU) rRNA

**the oldest functional part, conserved throughout evolution**



## Last Universal Common Ancestor (LUCA)

**LSU of tRNA (E. coli) :**

**Legend:**

- Phase 6 (Red)
- Phase 5 (Orange)
- Phase 4 (Yellow)
- Phase 3 (Green)
- Phase 2 (Light Blue)
- Phase 1 (Dark Blue)

**LSU and SSU associated at the end of Phase 3**

**LSU and SSU evolved independently & uncorrelated through Phases 1-3**

**PTC** (Protein Transfer Center)

**Phase 3**

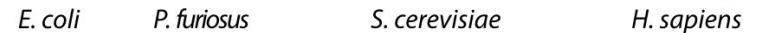
**Phase 4**

**Phase 5**

**Phase 6**

### LSU of tRNA (E. coli):

**Blue part is the oldest one**



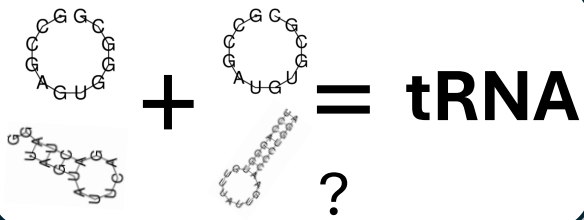
Molecular level chronology of the evolution of the large ribosomal subunit (LSU) rRNA. Each accretion step adds to previous rRNA but leaves the underlying **core unperturbed** (Anton S. Petrov et al., PNAS, 2015)

This symmetry (**SymR**) suggests that the ancient ribosome may have been **a dimer of identical or nearly identical RNA molecules**, later evolving into the asymmetrical modern ribosome with **PTC**.

# Peptidyl Transferase Center (PTC) Sequences

## the idea

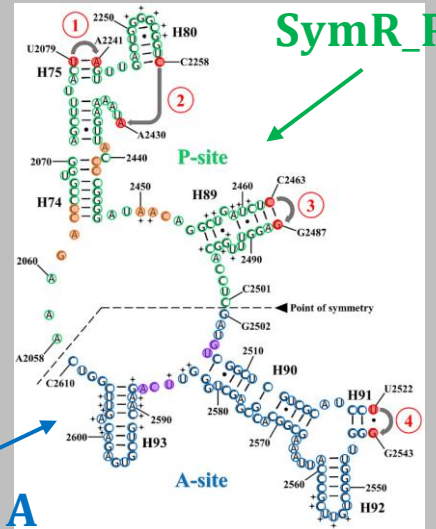
the dimerization of two similar RNA structures



“The peptidyl transferase center (PTC) evolved from a primitive system in the RNA world comprising tRNA-like molecules formed by **duplication of minihelix-like small RNA**”

Tamura, J. Biosci, 2011

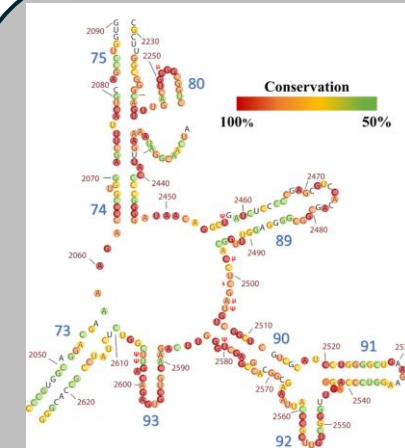
## pseudosymmetrical region



$$\text{SymR}_{\text{PA}} = \text{SymR}_{\text{P}} + \text{SymR}_{\text{A}}$$

Secondary structure of the pseudosymmetrical region (**SymR**; Agmon et al., 2005), derived from the LSU secondary structure of *Thermus thermophilus* (Petrov et al., 2013). (Madhan R. Tirumalai et al., 2021)

## PTC



PTC2 = red  
PTC3 = PTC2 + orange  
PTC4 = PTC3 + yellow  
PTC5 = PTC4 + green

**Nucleotide CONSERVATION level:**

Red circles: 100% conservation (78 nt).  
Orange circles: 90 to 99.9% conservation (68 nt)  
Yellow circles: 70 to 89.9% (52nt)  
Green circles: 50 to 69.9% conservation (49nt)  
Black letters: less than 50% conservation (35nt)

(Bernier et al.; Faraday Discuss, 2014)  
(Madhan R. Tirumalai et al., 2021)

**SymR\_P** is older than **SymR\_PA**  
**PTC2** is older than **PTC3**, **PTC4**, **PTC5**