

Co-transcriptional quantification of RNA synthesis, RNA folding and RNA-protein interactions by "Systems NMR" approach.

Yaroslav Nikolaev (✉ yaroslav.v.nikolaev@gmail.com)

ETH Zurich <https://orcid.org/0000-0002-1479-7474>

Nina Ripin

<https://orcid.org/0000-0001-6099-6506>

Martin Soste

Paola Picotti

<https://orcid.org/0000-0002-4109-3552>

Dagmar Iber

<https://orcid.org/0000-0001-8051-1035>

Frédéric H.-T. Allain (✉ allain@mol.biol.ethz.ch)

ETH Zurich <https://orcid.org/0000-0002-2131-6237>

Keywords: Systems Biology, NMR, biochemical networks, transcription, RNA folding, interactions, phase separation, liquid-liquid-phase-separation, NMR spectroscopy, RNA-binding proteins

Posted Date: August 22nd, 2019

DOI: <https://doi.org/10.21203/rs.2.9160/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

Abstract

The protocol describes how to setup and analyse observation of a co-transcriptional RNA folding network by Systems NMR approach. While most experimental approaches can monitor only a single molecule class or reaction type at a time, Systems NMR permits single-sample dynamic quantification of entire "heterotypic" networks – involving different reaction and molecule types. It thus provides a deeper systems-level understanding of biological network dynamics by combining the dynamic resolution of biochemical assays and the multiplexing ability of "omics".

This particular protocol describes the reconstruction of an 8-reaction co-transcriptional network - with simultaneous monitoring of RNA, metabolite, and proteins in a single sample at the same time. From reactions side, the protocol simultaneously quantifies RNA transcription, RNA folding and RNA-protein interactions (observed both from RNA and from protein side) and few other auxiliary reactions. In addition to fundamental analyses of reaction constants under different conditions, the current applications of this particular reconstruction are: (1) map RNA-binding interfaces on proteins without having to purify/order the RNA; (2) monitor co-transcriptional RNA folding perturbations by proteins and small molecules; (3) monitor RNA-transcription-driven protein phase-separation with the possibility to observe multiple proteins at once, each with residue-level resolution. Not counting the protein and RNA template preparation times, the NMR measurement and data analysis parts take about 1 day each.

This protocol accompanies Nikolaev et al, Nature Methods, 2019 (doi:XXX). The most up-to-date version of the protocol (including example code and data) is available at: github.com/systemsnmr/ivtnmr

Introduction

Reagents

DNA/RNA template preparation:

- Standard DNA cloning, bacterial expression and DNA purification reagents
- Midi/Maxi kit for DNA purification allowing to purify ≥ 250 ug plasmid DNA
- NEB buffer 3.1
- BsaI enzyme for plasmid linearization

In-vitro-transcription in NMR tube

- 5mm (or 3mm) NMR tubes (TA type is ok)
- 20x transcription buffer (TTD77): 800 mM Tris-HCl, pH 7.7, 0.2% Triton X-100, 100 mM dithiothreitol (DTT).
- 0.5mM DSS (4,4-dimethyl-4-silapentane-1-sulfonic acid) solution in D2O
- 80 mM nucleotide triphosphate (NTP) solutions (for each NTP, Applichem, ATP: A1348.0005; GTP: A1803.0001; CTP: A2145.0500; UTP: A2237.0001) - adjusted to pH 7.7-8.
- 1M MgCl₂
- 100 U/ml Inorganic Pyrophosphatase from baker's yeast (Sigma)
- ~70 uM (7 g/l 100 kDa) - 250x T7 RNA Polymerase (usually prepared in-house)
- [If protein observations are planned]: RNA-binding protein of interest

Equipment

- Standard DNA cloning, bacterial expression and DNA purification equipment
- NMR spectrometer with 1H + 31P cryoprobe (e.g. CPQCI), and ≥ 500 MHz 1H frequency (~12 Tesla)
- Computer with TopSpin and Matlab software

Procedure

The most up-to-date version of the protocol (including example code and data) is available at: github.com/systemsnmr/ivtnmr

(A) Design and preparation of RNA transcription template(s), and protein (~8 days)

1. RNA sequence design (~0.5 day)

Because substantial fraction (30-70%) of transcription products are short abortive RNA products, an ~8nt-long 5'-overhang nucleotide sequence needs to be prepended to the main RNA, to minimize interference of these short abortive RNAs with specific protein-RNA interactions and RNA folding. This sequence is designed algorithmically and can be used as a separate control to identify specific RNA effects from the other network perturbations.

Having decided on the primary RNA sequence, run $samp \leq_{pre} \frac{P}{a} b$ or $ts_n - backfold \in g_X X.m$ script, following instructions/examples in the header of the script for 5'-overhang generation (e.g. avoiding formation of dimers, excluding GG/purine pairs, etc. The script will produce a range of possible 5'-overhang variants, from which one or several can be selected for experimental tests. Main criteria for downstream experimental selection are:

- Judging from denaturing / non-denaturing gel-electrophoresis:

- homogeneity of transcription product

- yield of transcription product

- If observation of imino signals from the folded RNA is important, and you're choosing between several 5'-overhang variants: run 1D1H NMR of the transcription mixture, and choose the sequence which produces the least number of "background" signals in the imino region of the spectrum.

2. Preparation of DNA templates (~7 days)

To ensure maximum homogeneity of synthesized RNA, we do the transcription from a plasmid DNA template (not from synthetic oligo-nucleotide templates). To increase efficiency, multiple templates can be prepared in parallel.

2A. Clone the template DNA under T7 RNA Pol promoter (~1 day cloning, ~1 day wait for colonies, ~2 days miniprep + sequencing)

We are using pTX1 vector system [Michel E. et. al, 2018](https://doi.org/10.1007/978-1-4939-7634-8_11). The $samp \leq_{pre} \frac{P}{D} tx1 . rs. xls$ provides automated template for design of cloning primers (see instructions in the file).

2B. Purify DNA template (~1.5 days)

- Use a kit Nucleobond Xtra Maxi (Macherey-Nagel) or an equivalent from Qiagen, etc. This yields per column ≥ 250 -350 ug pTX1 plasmid DNA template, enough for ~10 transcription-NMR reactions (450ul each, at 33 nM DNA template).

Tips for the procedure:

- Per one Maxi column grow ~300ml culture till $OD_{600} \approx 3$ ($OD * V = 900$).

- Can try to saturate the Maxi columns more, to reach 500-1000 ug plasmid yield per column: then grow e.g. 600ml of $OD_{600} \approx 3$ culture, and increase the volumes of buffer used for lysis & etc.

- To get max yields: start the culture use a fresh (<7 days) colony, make overday preculture in 10ml until $OD_{600} \approx 1$, dilute 1:100 in the final ~300 ml medium for overnight growth. (or store the fresh $OD_{600} \approx 1$ preculture at 4°C and start main culture in the morning).

- Use baffled flasks, this allows to reach $OD_{600} \approx 3$, giving 4g cells / L culture.

- Use 2x antibiotic amount than usual.

- Elution from column: pre-heat elution buffer + make elution twice (reapplying same eluate).

- After isopropanol precipitation centrifuge 1h @ 4000xg, 4°C.
- After isopropanol removal, wash the pellet not once, but 3 times with 70% ethanol, each time detaching DNA pellet from the walls by inverting the tube, and then centrifugation. This removes salts which reduce the yield &/or quality of transcription.
- Dry the pellet on air for 4-6 hours (overnight also ok).
- Dissolve DNA in H₂O to ~0.41 μM concentration (650 μg/ml for 2.4kb plasmid) - this allows to have sufficiently concentrated DNA in the end (after linearization), to not dilute the final NMR sample too much.
- To ensure homogenous dissolution of DNA pellet after drying, perform 3 cycles of: [vortex, incubate 10-20 minutes at 50°C, vortex, freeze at -20°C, thaw]. If pellet did not detach from the wall after first vortex, use the pipet to detach it.

2C. Linearize the DNA template (~0.5 day).

To maximize the yield and homogeneity of the main RNA product, it is important to thoroughly linearize the template. For the pTX1 plasmid we use BsaI enzyme (NOT the BsaI-HF version!), at 1.5U/ug DNA, in NEB buffer 3.1, for 13-15 hours at 50°C, followed by 1h at 65°C inactivation. Check the plasmid cleavage on 1% agarose gel.

2D. Optimize the MgCl₂ concentration for transcription (~1 day)

Run small-scale transcriptions (~25 μl is enough) and analyze them using denaturing PAGE with urea. In our hands for pure DNA template optimal MgCl₂ range was always close to 1:1 ratio with NTPs (e.g. 20-24mM MgCl₂ when using 20mM NTPs).

3. Buffer exchange for protein of interest (~0.5 day)

To match the starting Tris-Triton-DTT (TTD) transcription buffer, transfer your protein into: 40 mM Tris-HCl, 0.01% Triton-X100, 5 mM Dithiothreitol, pH 7.7. From our experience many proteins are stable in this buffer. To stabilize the protein more, one can potentially add:

- 50 mM L-Arg/L-Glu (this leads to increased broadening of imino signals)
- up to ~75 mM NaCl / KCl (salt decreases transcription efficiency)
- ~1 mM NTPs - these seem capable to mask RNA-binding interfaces, reducing self-aggregation of protein through those. Added NTPs need to be taken into account in the total NTP concentration – for addition of MgCl₂, and for setting the initial concentrations during ODE network modeling.

(B) Setup of an in-vitro-transcription-NMR reaction (30-60 minutes)

1. Sample prep

- Sample preparation template: $samp \leq_{pre} \frac{P}{i} vtnmr. xslx$

- Below setup is for 450 ul samples in 5mm TA NMR tubes. One can also use 250-300ul samples in shigemi tubes and/or 150ul in 3mm tubes (but for 3mm tubes the acquisition times may need to be adjusted).

- To increase observability of imino signals we:

- Run reactions at 30°C instead of 37°C.
- Exclude commonly added spermidine from the transcription buffer because it leads to significant broadening of imino signals.
- Use transcription buffer with pH 7.7. T7 RNA Pol gives better yields at pH 8.1, but iminos are observed better at lower pH. During typical transcription reaction (starting with 20mM NTPs), due to release of free xPO₄, pH goes down by ~0.2 pH units. The pH change can be monitored from the shift of 31P xPO₄ signal and/or shift of 1H Tris signal.

- In case you need larger fraction of sample volume to add diluted protein or other additives, the NTPs, MgCl₂ and DNA template components can be pre-lyophilized together. This seems to not affect the yield or homogeneity of RNA product.

- T7 RNA polymerase is added only after time0 reference spectra are recorded.

2. NMR setup

Notes:

!!!!!!!

WARNING: the TopSpin Python scripts and pulseprograms are provided here "as is" - merely as a guideline for automated setup. These were not thoroughly tested on different spectrometers and may contain bugs and incompatibilities with TopSpin and spectrometer console versions!

!!!!!!!

- We use an automated script $nm\frac{r}{p}\frac{y}{\in v}02$. *yn.py* which makes tuning-matching, shimming, pulse calibrations and experiment setup.

- The measured experiments are configured by this script based on the template experiments which are placed into an empty IVTNMR_template - template dataset (example in $nm\frac{r}{I}VTNMR_t\text{emplate}$):

- 1D1H (expno 12, pp $nm\frac{r}{p}\frac{p}{z}g$ - *wg001*) - full 1H spectrum - for DSS calibration, pH check from Tris position, and potential quantification of individual NTPs (A/U/G/C) consumption rates (all four nucleotides have some specific signals in the aromatic region of the spectrum).

- 2D 1H1H TOCSY (expno 13, pp $nm\frac{r}{p}\frac{p}{s}$ → *csy003*) - (not required for the setup described in the paper) - allows to observe RNA with higher resolution on certain signals, e.g. U/C H5-H6 correlations.

- 1D31P (expno 14, pp $nm\frac{r}{p}\frac{p}{z}gig002$) - to observe 31P-containing molecules

- 2DHN-sofast-hmqc (expno 15, pp $nm\frac{r}{p}\frac{p}{h}mqc01$) - to observe protein

- 1D1H-sofast (expno 16, pp $nm\frac{r}{p}\frac{p}{z}g$ - *sof * 006*) - to observe RNA imino region with increased sensitivity

Step-by-step procedure:

1. Create new dataset, including main information in its name, for simplification of automated analysis:

$YYMMDD_I NXXX_R NANAME_P ROTENAME_T EMPERATURE_M AGNET$ (e.g. $180914_I N71 b_S MN1_c oNUP1_{303} K_{600}$). Here, INXXX (IN71b) - is the experiment ID, implying same ID for the same RNA and protein combination. Replicate experiments are denoted with lowercase letter (IN71a, IN71b, ..) - so its easy to find similar experiments programmatically from shell, Matlab, etc.

2. Make temperature calibration (we use expno 1 consistently).

3. Insert the sample into magnet. Wait 1 min for temperature equilibration. Run $nm\frac{r}{p}\frac{y}{\in v}02$. *yn.py* script for the setup of series. Check and test this script carefully before doing real runs.

4. After the automation script has finished:

- In 2DHN spectrum - set the SW / carrier / TD / etc parameters to the values optimal for your protein

- Check the 90-degree 31P pulse value (script runs *paropt* procedure, storing results in expno 4 999) - and enter it into 1D31P experiment.

- (If expect drift of tune-matching system): check the tune-match of channels of interest.

3. Process 1D31P:

- $qu\mu < i5000 - 5500$

- example of processing command (check only SR corrections, in our setup phase seems determined robustly enough by *apks* routine):

SR - 139; si64k; wdwEM; lb2; efp; apks; |f114; |f2 - 26; |g|5; |n|

4. Integrate: TopSpin \int *ser* command, pointing to the \int *egr_datasets31P.txt* as the list of spectra to integrate (it should be created by *s* or *t.py* in the dataset folder). Use default settings for integration - the calibration of integrals to internal NTP signal happens later during analysis.

- Integration regions we use:

3.4714285714285715 1.5571428571428572 – PO4

0.42857142857142855 -2.676190476190476 – PPI

-4.3238095238095235 -5.352380952380952 – RNA

-5.352380952380952 -5.390476190476191 – gammaNTP

-5.390476190476191 -5.804761904761905 – RNA-5'gamma?

-9.101442036015632 -9.806884238970339 – betaNDP

-9.80952380952381 -10.433333333333334 – alphaNDP

-10.433333333333334 -11.219047619047618 – alphaNTP

-17.60952380952381 -19.052380952380954 – RNA-5'beta?

-19.052380952380954 -19.714285714285715 – betaNTP

4. Process 2DHN:

- $qu\mu < i4000 - 4500$

- example of processing command (check phase, SR corrections, STSI/STSR for the dimensions of the final spectrum, etc):

2sr - 26.41; 1sr - 2.676; 2phc0 - 264.203; 2phc10; 2si8k; 1si512; 2STSI2214; 2STSR1512; 1stsi0; 1stsr0; 2|f111; 2|f25.5; 1|f1200; 1|f280; |g|5.

5. Pick peaks in 2DHN spectra: create a project and a peaklist in a Cara repository, then pick peaks of interest and trace their positions across time series using peak aliasing in Cara Monoscope. The *s* or *t.py* script creates a cara repository with linked spectra in the dataset folder. Downstream analysis scripts read the peaks information from cara repositories at the moment. If analyzing multiple datasets at once, it is, however, convenient instead of having multiple Cara files, to have one Cara file with each dataset represented by own project and peaklist. Useful shortcuts for peak tracing in Monoscope:

- $Ctrl + \frac{1}{C}trl + 2$ - move between spectra

- $m\frac{p}{m}a$ - move peak or its alias

- *gp* - go to peak

- *gs* - go to spectrum (if not argument provided will switch to the first spectrum in series)

(E) Data analysis and model fitting (≈ 0.5 day)

Unless stated otherwise, most below procedures are done in Matlab (.m files).

1. Generate data structure for ODE model fit.

Run `analysis/ODE/v01/mod_01_el_and_dat_01_ivtnmr.m` - this script will read 31P integrals from above 31P integration file, and 2DHN chemical shifts from cara repository, and then create the data structure used by the ODE model fitting procedure in `ODE/v01/mod_01_el_and_dat_01_ivtnmr` folder. It will also generate "ivtnmr" data object in `ODE/v01/model_and_data/data_ivtnmr_full/` - which contains main information about IVTNMR experiment in one Matlab structure (names, number of time-points, integrals, chemical shifts, etc). The code in v01 folder is largely self-contained - so its convenient to just duplicate and rename it (e.g. v02_test) to keep track of different versions when you're making adjustments to data analysis / model structure / etc.

2. Fit the ODE model.

Run `analysis/ODE/v01/a_fit_multi.m`, which can fit multiple IVTNMR datasets sequentially. The procedure reads the ODE model from `model_and_data/xIVT.mdf` file and fits it into experimental data from `model_and_data/data_for_fit` generated at the previous step. You may need to adjust the `ma` length in the MDF file to suit your RNA length - it is included in the <RHS> section (e.g. for 28-nucleotide RNA: `+f1*(27/28) +f2 -f3`). The `mdf` file is generated by Matlab `convertBNGL_to_MDF.m` script based on the `xIVT.bngl` - Rule-Based-Model defined in BioNetGen language (mode details below). The fitting script will generate PDF figures of the fit summary (in `ODE/v01/model_and_data/figure_images`) and will export fitting results into `ODE/v01/model_and_data/a_fit_multi.mat` which can be read and analyzed further.

3. Visualize ODE fit results.

Run `analysis/ODE/v01/a_visualize.m` - this will automatically read fitting results from `ODE/v01/model_and_data/a_fit_multi.mat` and visualize the fitted constants, replicates and errors.

4. Analyze and visualize time-resolved imino linewidths.

Run `analysis/LW/a_fit_LW.m` to fit and visualize imino linewidths. This fitting assumes a single peak with lorentzian lineshape, and thus works best for well-resolved signals. Decrease of pH during reaction (usually around ~ 0.2 pH units) might lead to slight systematic change of linewidth with reaction time. Also, a "growing" baseline in a crowded region may lead to apparent narrowing of linewidth with reaction time, because the current fitting routine assumes baseline fixed at zero intensity. This effect can be factored out, if appropriate baseline correction is found. In the single-lorentzian fit, the intrinsic broadening of one imino signal cannot be distinguished from broadening caused by overlap with neighboring signals, which has to be considered when interpreting the data.

(F) Analysis of parameter (reaction constant) uncertainty

-**Replicates.** Based on our current experience, the most sensible / realistic parameter uncertainties are obtained by running ≥ 2 experimental replicates - ideally using different batches of DNA template and/or protein. Basic code for calculation of such uncertainties is included in `analysis/ODE/v01/a_visualize.m` and `analysis/LW/a_visualize_LW.m`.

-**Bootstrap.** Uncertainties obtained from bootstrap analysis (resampling of the full data vector with replacement prior to the fitting), in our setup yield too narrow confidence intervals - most likely due to very high number of recorded data points.

-**FIM.** The standard Fisher Information Matrix / Cramer Rao bound also often yields unrealistically narrow or unrealistically large uncertainties - most likely due to the lack of reasonable error estimates on the individual NMR integral / chemical shift data points.

(G) Adjusting ODE model

1. In [RuleBender](<http://visualizlab.org/rulebender/>) software (or just text editor): **edit the BioNetGen network** definition (`ODE/v01/model_and_data/xIVT.bngl`). For downstream fitting it is important to match the number and order of defined observables to the data vectors you actually provide for model fitting.

2. In shell/terminal: ****recompile the model****.

```
cd ODE/v01/model_and_data
../../soft/BioNetGen217/Perl2/BNG2.pl xIVT.bngl
```

3. In Matlab: ****convert the model to MDF format**** for fitting.

```
cd ODE/v01/model_and_data
convertBNGL_to_MDF('xIVT.bngl')
```

4. In Matlab: ****adjust the fitting routine**** v01/a_fit.m to match the observables / data / constants used in the new network model. Key things which need to be checked:

- Order/number of data vectors used for fitting (needs to match the observables in mdf file) - check obs_data_to_keep variable;
- If added new observables/reactions: check if model initial params (NDP, Prot(s) conc, ..) are set correctly
- If added new constants/params: check their init values in the a_fit.m, file and/or exclude from optimization. And potentially add their normalization and saving into the output of model fit results.

****Note****: In the used here implementation BioNetGen cannot provide arbitrary rate laws and algebraic constants for scaling. If want such things included - need to specify them in the final *.mdf file used for fitting - similar to what we do for RNA length definition. Alternatively one can switch to e.g. [PottersWheel] (<https://potterswheel.de/>) software which provides more advanced, up-to-date and fully integrated environment for model definition and optimization, including model optimization in log-space.

Disclaimer: Limitations of Liability for the code

The code, pulseprograms, and especially TopSpin Python scripts, in this repository are provided "as is" - merely as a guideline for automated setup. These were not thoroughly tested on different spectrometers and may contain bugs and incompatibilities with TopSpin versions. Authors assume no responsibility, and shall not be liable to you or to any third party for any direct, indirect, special, consequential, indirect or incidental losses, damages, or expenses, directly or indirectly relating to the use or misuse of the code and pulseprograms provided here.

Troubleshooting

Indicated directly in corresponding sections.

Time Taken

Indicated directly in corresponding sections.

Anticipated Results

Key results from the network reconstruction:

- Overall pattern of RNA iminos - reporting on the structural fold of the RNA.
- Time-resolved LineWidths of RNA imino signals - quantitatively reporting on structural in RNA during transcription reaction. Linewidths are convertible to ΔG - free energy of folding, when the intrinsic and base-flipping rates of imino exchange ($k_{ex,intrinsic}$ and $k_{ex,base-flipping}$) are known.

Processing math: 64%

- Time-resolved and RNA-concentration resolved chemical shift and signal intensity perturbations for all observable residues in the protein of interest - reporting on potential interaction sites / structural changes in the protein.
- K_D of the protein-RNA interaction (if protein becomes saturated during reaction, and/or if chemical shifts of protein residues in RNA-saturated state are known).
- k_{cat} rate of RNA transcription.
- Individual rates of ATP, GTP, UTP, CTP consumption (currently not yet implemented in the automated analysis).

References

Acknowledgements

We thank J. Vollmer and G. Fengos for the help with network modeling. We acknowledge G. Wider and all members of the ETH BNSP platform for excellent maintenance of the NMR infrastructure. We thank all members of the Allain Lab, in particular F. Damberger, and the Parpan retreat participants for helpful discussions. This work was supported by the Promedica Stiftung, Chur (Grant 1300/M to Y.N.), Novartis Foundation and Krebsliga Zurich (Y.N.), NCCR RNA and Disease by the Swiss National Science Foundation (F.A. and N.R.).

Figures

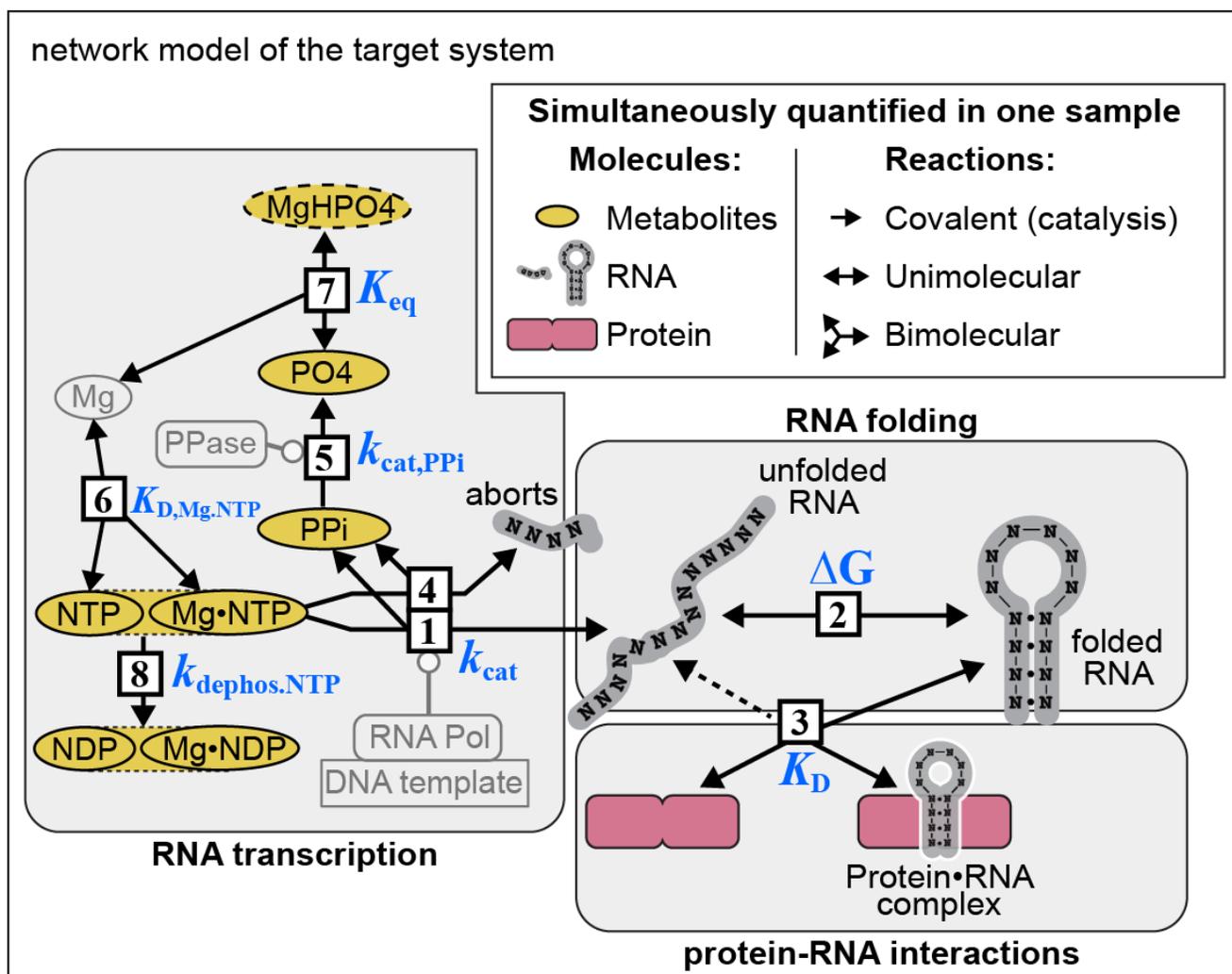


Figure 1

Network model of the co-transcriptional system monitored/quantified in this Systems NMR assay.

Supplementary Files

Supplementary files associated with this preprint. Click to download.

Processing math: 64%

- supplement2.zip
- supplement3.mp4
- supplement3.png