

Estimating RNA dynamics using one time point for one sample in a single-pulse metabolic labeling experiment

Micha Hersch, Adriano Biasini, Ana C. Marques, Sven Bergmann

July 17, 2021

Abstract

Over the past decade, experimental procedures such as metabolic labeling for determining RNA turnover rates at the transcriptome-wide scale have been widely adopted and are now turning to single cell measurements. Several computational methods to estimate RNA processing and degradation rates from such experiments have been suggested, but they all require several RNA sequencing samples. Here we present a method that can estimate RNA synthesis, processing and degradation rates from a single sample. Our method is computationally efficient and outputs rates that correlate well with previously published data sets. Using it on a single sample, we were able to reproduce the observation that dynamic biological processes tend to involve genes with higher metabolic rates, while stable processes involve genes with lower rates. This supports the hypothesis that cells control not only the mRNA steady-state abundance, but also its responsiveness, i.e., how fast steady-state is reached. In addition to saving experimental work and computational time, having a sample-based rate estimation has several advantages. It does not require an error-prone normalization across samples and enables the use of replicates to estimate uncertainty and perform quality control. Finally the method and theoretical results described here are general enough to be useful in other contexts such as nucleotide conversion methods and single cell metabolic labeling experiments.

1 Introduction

Since the advent of molecular biology, a consensus has emerged that the regulation of gene expression underlies most biological processes including development, disease and adaptation [1, 2, 3]. While gene expression regulation has mostly been associated with activating the production of RNA (e.g. through transcription factors), it has become apparent that the regulation of RNA splicing and RNA stability also plays an important role in determining the expression level of a gene [4, 5]. Taking advantage of high throughput RNA quantification

protocols, methods designed to distinguish the effects of RNA synthesis, processing and degradation at the transcriptome-wide level have been developed. Among them, RNA metabolic labeling techniques relying on chemically modified ribonucleotides such as 4-thiouridine (4sU) and 5'-Bromouridine (BrU) have been widely adopted, as their impact on cellular function is minimal [6]. Briefly, incubating cells with modified ribonucleotides for a limited period of time (referred to as the pulse), and their concomitant incorporation in newly synthesized transcripts, allows distinguishing newly transcribed from preexisting RNA, which can be biochemically separated and quantified. This quantification, initially performed using microarray [7] and later using RNA-seq [8, 9], can then be used to estimate RNA decay rates. More recently, methods that rely on nucleotide conversion have been used to the same effect, with the advantage of circumventing the cumbersome biochemical enrichment and separation step.

In the last decade, several methods to estimate RNA dynamics from metabolic labeling experiment data have been developed [10, 11, 12] (see [13] for a review). Typically, labeled transcript abundances are fitted to an exponential function approaching to steady-state equilibrium (during or after the pulse), from which the RNA half-life can be estimated [14, 15, 16]. This requires time-course experiments in order to have enough points for fitting, as well as a way to normalize RNA concentrations across samples, either using spike-ins [17], or using internal controls such as intron concentrations [18]. The INSPEcT method [19] goes beyond first order dynamics and takes into account the RNA processing rates, which are estimated along with the degradation and synthesis rates. This method works by first estimating rates for individual samples by assuming no degradation during the pulse and then uses those estimates as a starting point for fitting model of rate evolution for all the rates of all samples. Those methods rely, for each sample, on a the separate quantification of labeled RNA on one hand and of total (mixed labeled and unlabeled) and/or unlabeled (or pre-existing) RNA on the other hand. In its later version, INSPEcT was extended to estimate rates without labeling the sample [20].

In this work, we build on the INSPEcT approach and derive an exact solution (when it exists) for the initial rate estimates without making the assumption of no labeled transcript degradation. This is achieved by considering the intron to exon ratio for each transcript in both the labeled and unlabeled RNA pools, thus allowing to bypass the need for normalization across those two samples. We can thus infer synthesis, processing and degradation rates from a single sample and time point. Those rates can be used as such, allowing to reduce the experimental load and costs and compare rates across samples and time points. But they can also be used, as in INSPEcT, as initial estimates for multiple sample-based rate estimation. Applying our method to our own experimental data and using a single sample and time point, we obtain synthesis and processing rates that are well correlated with the ones obtained using INSPEcT first guess. The degradation rates, on the other hand, correlate poorly across the two methods, but those computed with our method correlate better than INSPEcT with pre-

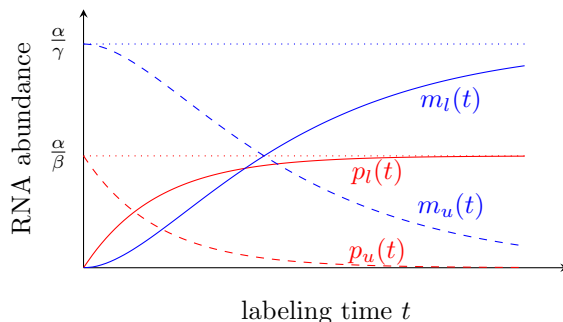


Figure 1: Evolution of unlabeled and labeled, premature and mature RNA during labeling according to the Zeisel model. Dotted horizontal lines correspond to steady-state levels, dashed lines correspond to the unlabeled RNA and solid lines to labeled RNA. Processing and degradation rates can be estimated from the ratios of the two dashed lines and of the two solid lines at a single time point.

viously published mRNA degradation rates obtained with three replicates and seven time points [21]. Because it can be reduced to numerically solving an equation with a single unknown on a bounded domain, it is also faster than INSPEcT. Moreover, our results are consistent with an adapted gene-specific mRNA responsiveness and co-transcriptional mRNA processing [22].

2 Method

2.1 Overview

This paragraph summarizes the general strategy of the method, with references to relevant equations indicated in parentheses. We use the Zeisel model of RNA dynamics [23] to model both the unlabeled and the labeled RNA (1, 2). Using the standard procedure for solving systems of linear differential equations, we find its general solution and its free parameters by setting the initial conditions for both the unlabeled (or pre-existing) and the labeled RNA (3-6), as illustrated in Fig 1. We can then express, for a given gene, the ratios for both unlabeled and labeled RNA of intron to exon expression level as functions of the processing and degradation rate of that gene (8,9). These two ratios are independent from the RNA synthesis rate. Using the intron to exon ratios as observables, we are left with two non-linear equations and two unknowns, namely the processing and degradation rates. These equations are then reparametrised with dimensionless parameters and reduced to a single non-linear equation with one unknown (12). This resulting equation is only defined on a bounded domain (13). Our rates can thus be inferred by numerically solving that equation on a bounded domain, which is very fast. In addition, we prove in Appendix C that this equation, under

certain conditions, has a single solution (but in general it can also have two or no solution).

2.2 Model

Like previous work [19], we use the Zeisel model of RNA synthesis, processing and degradation [23].

$$\dot{p} = \alpha - \beta p \quad (1)$$

$$\dot{m} = \beta p - \gamma m, \quad (2)$$

where p is the premature RNA, m the mature RNA, and α , β , γ are RNA the synthesis, processing and degradation rates. This model can be solved analytically (see appendix A). In particular, enforcing the boundary conditions corresponding to the unlabeled RNA, namely that it is at steady-state when the pulse starts ($t = 0$) and that subsequently no more pre-mature RNA is produced, results in

$$p_u(t) = \frac{\alpha}{\beta} \exp(-\beta t) \quad (3)$$

$$m_u(t) = \frac{\alpha}{\gamma - \beta} \exp(-\beta t) - \frac{\beta \alpha}{\gamma(\gamma - \beta)} \exp(-\gamma t), \quad (4)$$

where the u subscript indicates that this corresponds to the unlabeled RNA pool.

Enforcing boundary conditions corresponding to the labeled RNA, namely that it is not (yet) expressed at $t = 0$ leads to

$$p_l(t) = \frac{\alpha}{\beta} (1 - \exp(-\beta t)) \quad (5)$$

$$m_l(t) = \frac{\alpha}{\gamma} \left(1 + \frac{\beta}{(\gamma - \beta)} \exp(-\gamma t) \right) - \frac{\alpha}{\gamma - \beta} \exp(-\beta t) \quad (6)$$

where the l subscript indicates that this corresponds to the labeled RNA pool.

2.3 Inferring synthesis, processing and degradation rates

We consider that the exonic RNA abundance χ corresponds to the premature and mature RNA, while the intronic RNA abundance ι correspond to the premature RNA only. Furthermore, we assume that χ and ι are suitably normalised for exonic and intronic length so that they are proportional to the number of transcripts. We can then compute:

$$\frac{\iota}{\chi} = \frac{p(T)}{p(T) + m(T)}, \quad (7)$$

where T is the duration of the labeling.

In the case of unlabeled fraction, we have

$$\begin{aligned}
 \frac{\iota_u}{\chi_u} &= \frac{p_u(T)}{p_u(T) + m_u(T)} \\
 &= \frac{E_\beta}{\beta\left(\left(\frac{1}{\beta} + \frac{1}{\gamma-\beta}\right)E_\beta - \frac{\beta}{\gamma(\gamma-\beta)}E_\gamma\right)} \\
 &= \frac{E_\beta}{\frac{\gamma}{\gamma-\beta}E_\beta - \frac{\beta^2}{\gamma(\gamma-\beta)}E_\gamma} \\
 &= \frac{(\gamma - \beta)E_\beta}{\gamma E_\beta - \frac{\beta^2}{\gamma}E_\gamma} \\
 &= \frac{\gamma(\gamma - \beta)E_\beta}{\gamma^2 E_\beta - \beta^2 E_\gamma}
 \end{aligned} \tag{8}$$

where we define $E_\beta = \exp(-\beta T)$ and $E_\gamma = \exp(-\gamma T)$ as abbreviations.

For the labeled fraction, we have

$$\begin{aligned}
 \frac{\iota_l}{\chi_l} &= \frac{p_l(T)}{p_l(T) + m_l(T)} \\
 &= \frac{(1 - E_\beta)}{(1 - E_\beta) - \frac{\beta}{\gamma-\beta}E_\beta + \frac{\beta}{\gamma}\left(1 + \frac{\beta}{\gamma-\beta}E_\gamma\right)} \\
 &= \frac{(1 - E_\beta)}{\frac{\gamma+\beta}{\gamma} - \frac{\gamma}{\gamma-\beta}E_\beta + \frac{\beta^2}{\gamma(\gamma-\beta)}E_\gamma} \\
 &= \frac{\gamma(\gamma - \beta)(1 - E_\beta)}{\gamma^2 - \beta^2 + \beta^2 E_\gamma - \gamma^2 E_\beta} \\
 &= \frac{\gamma(\gamma - \beta)(1 - E_\beta)}{\gamma^2(1 - E_\beta) - \beta^2(1 - E_\gamma)}.
 \end{aligned} \tag{9}$$

We notice that this last expression is of the same form as the one for the unlabeled fraction (8), but replacing exponentials by their complement to one. Importantly these two fractions do not depend on α , which (unlike [20]) allows our method to estimate processing and degradation rates independently from the synthesis rate.

Denoting $a = \frac{\iota_u}{\chi_u}$ and $b = \frac{\iota_l}{\chi_l}$ as the observable unlabeled and labeled fractions of intron abundance, we are left with a system of two equations and two unknowns β and γ , which we now set out to solve. First, we reparametrise our system with $\beta = k\gamma$ and define $E_{k\gamma} = E_\beta = \exp(-k\gamma T)$ leading to

$$a = \frac{(1 - k)E_{k\gamma}}{E_{k\gamma} - k^2 E_\gamma} \tag{10}$$

$$b = \frac{(1 - k)(1 - E_{k\gamma})}{(1 - E_{k\gamma}) - k^2(1 - E_\gamma)}. \tag{11}$$

It is shown in Appendix B that system of equations can be simplified to the following equation in k :

$$\frac{k}{k-1} \log\left(\frac{k+a-1}{k^2 a}\right) - \log\left(\frac{(b-a)}{a(bk+b-1)}\right) = 0, \quad (12)$$

with the following domain of definition \mathcal{D} for k :

$$\max\left(\frac{1}{b} - 1, 1 - a\right) < k < \frac{1}{a} - 1. \quad (13)$$

The above equation does not explicitly depend on T can be solved numerically on \mathcal{D} . In practice a and b are approximated by r_u and r_l , defined as the length-normalized intronic to exonic read count ratio (or TPM ratio) for the unlabeled and for the labeled sampled respectively.

We further prove in Appendix C that for $b > \frac{1}{2-a}$, Eq. (12) has a single solution in the domain given by (13), which can be found very efficiently. This enables the estimation of the processing and degradation rates for a single sample. Moreover, since the reduced equation is independent from T , uncertainty on its true value does not affect the relative values of the resulting rates. Hence replicates can be used to assess the reliability of the estimates and time courses allow to test whether the rates are constant as assumed by the model.

If (12) does not have a solution, estimates can be obtained by minimizing (in log space) the squared Euclidean distance between the observed (i.e., r_u, r_l) and derived values of a and b :

$$f(k, \gamma T) = \left(\log(r_u) - \log\left(\frac{(1-k)}{\exp(-k\gamma T) - k^2 \exp((k-1) - \gamma T)}\right) \right)^2 + \left(\log(r_l) - \log\left(\frac{(1-k)(1 - \exp(-k\gamma T))}{(1 - \exp(-k\gamma T)) - k^2(1 - E_\gamma)}\right) \right)^2. \quad (14)$$

The ratios r_u, r_l must be smaller than one to make sense within our model and genes where this is not the case should be discarded. The log function is used to give exon and intron counts equal standing.

The above bivariate function can be reduced to a univariate function f^* using (36):

$$f^*(k) = f\left(k, \frac{1}{k} \log\left(\frac{r_l - r_u}{r_u(r_l k + r_l - 1)}\right)\right) \quad (15)$$

The processing and degradation rates are derived from k using (36) where a and b are again approximated by r_u and r_l respectively. Then the synthesis rate α can be easily obtained from (4), where m_u is approximated by χ_u (which is likely the most reliably measured specie):

$$\gamma = \frac{1}{kT} \log\left(\frac{r_l - r_u}{r_u(r_l k + r_l - 1)}\right) \quad \beta = k\gamma \quad \alpha = \frac{\gamma(\gamma - \beta)\chi_u}{\gamma E_\beta - \beta E_\gamma} \quad (16)$$

3 Results

3.1 Simulated data

In order to confirm that our method can be applied in principle, we evaluated our method on simulated data, where the data was generated using the exact model used to develop the method (see equations 3 and following). We chose not to simulate noise or biases in the data, as the aim of the simulation is only to validate the mathematical developments above and our implementation of the method. We generated 50000 random value for α , β , and γ ranging between $\exp(-5)$ and $\exp(5)$ and computed the corresponding values for ι and χ . We then computed r_u and r_l by taking the ratio. Estimates $\hat{\beta}$ and $\hat{\gamma}$ were then inferred by using r_u and r_l as an input to the method and compare the original β and γ .

Numerically solving equation (12), yielded either one or two solutions. The results for the unambiguous cases are shown in Fig 2, left. We see that in virtually all cases, the method yields accurate estimates of the processing and degradation rates. For a few points, the method is less accurate at the upper boundary of the parameter space, probably due limited floating point precision. Indeed if the labeling time is too long with respect to the metabolic rates, virtually all unlabeled RNA are degraded and the rates cannot be reliably estimated.

As we are considering single-sample estimates, it is possible to chart the observable space given by a and b and see when the method provides unambiguous results. Fig. 2, center, confirms that for $b > \frac{1}{2-a}$ the method provides a unique (and correct) solution as proven in appendix C. Below this line (displayed in green), the methods provides ambiguous results as two distinct set of values β and γ can account for the same value of a and b (in blue). It is also possible to visualize the trajectories of the observables a and b for various values of k , as depicted in Fig. 2, right. When $T = 0$, trajectories start from the top of the space at $(\frac{1}{1+k}, 1)$. When $k < 1$, as time passes the system moves down to $(a, b) \rightarrow (1 - k, \frac{1}{1+k})$. For $k \geq 1$, trajectories move to $(0, \frac{1}{1+k})$. Note that this is the expected case, as the splicing of mRNA occurs in general faster than its degradation. Note that, in this case, trajectories cross below the green line, explaining why two solutions can be found for a single value of (a, b) . The speed at which the system follows these trajectories depends on γ .

3.2 Real data

In order to assess the performance of the method on real data, we applied our method on the 4sU labeling experiment described in [24]. Briefly, mouse embryonic stem cells were plated at a density of 40000 cells/cm² on gelatin-coated 10cm tissue culture plates and grown for approximately 14 hours. After addition of 4sU to the growth medium, cells were incubated at 37C for 10 minutes (10 minutes labeling pulse). RNA was then extracted and processed according to the protocol described in [25]. Reads that did not map to mouse ribosomal

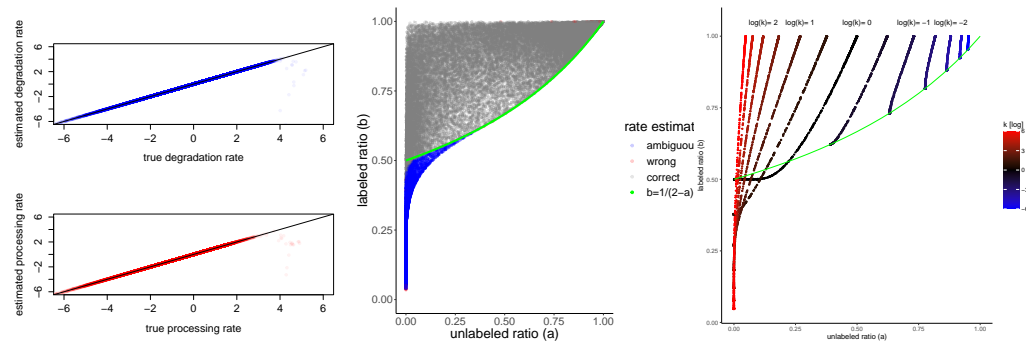


Figure 2: Simulated data. Left: the method correctly estimates processing and degradation rates. Points with ambiguous solutions are not shown. Some points corresponding to high rates cannot be estimated correctly as the system as already reached steady-state during the simulated "pulse". Center: the measurement space can be partitioned into ambiguous and unambiguous regions. The green line corresponds to $b = \frac{1}{2-a}$. Above that line, rates are correctly and unambiguously estimated. Boundary cases are sometimes wrongly estimated, probably due to numerical errors (red dots). Right: Trajectories in the phase space are solely determined by the k parameter. They start at time $T = 0$ at the top ($b = 1$) and go down. For $k < 1$ the trajectories (in blue) remain above the green line defined by $b = (2 - a)^{-1}$ and do not cross. For $k > 1$ (in red), they cross each other below the system follows the trajectory depends on the actual values of β and γ .

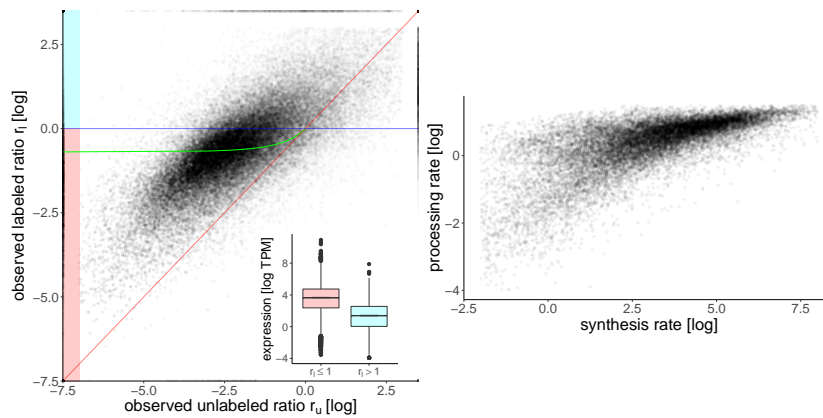


Figure 3: Real data. Left: Each point corresponds to a transcript with its transparency reflecting log expression value. Like in the previous figure, the green line is defined $y = (2 - x)^{-1}$. For transcripts lying between the abscissa (in blue) and the green line, estimates of processing and degradation rates can be obtained by solving (12). For transcripts lying between the diagonal (in red) and the green line, estimates can be obtained by minimizing (15). The observed ratios for the remaining transcripts are not coherent with the model and are discarded. These transcripts (above the blue line) are lowly expressed compared to the ones below the blue line (see inset). Right: RNA processing rates are highly correlated to the synthesis rates (72%), which is consistent with co-transcriptional RNA processing.

RNA sequences were aligned to intronic and exonic sequences using STAR V2.5 [26] and quantified using RSEM V1.1.17 [27], yielding intron and exon expression levels for unlabeled and labeled RNA.

For a single sample, the observable space represented in Fig 2 (center and right) is represented (in log coordinates) in Fig 3, left. We see that, while the points are centered on the expected region of the observable space, many transcripts lie below the diagonal or above the $r_l = 1$ (or $\log(r_l) = 0$) line (in blue), which is not compatible with our model. We observe that those incompatible transcripts lying above the $r_l > 1$ line are expressed at a much lower level than the transcripts lying below this line (see inset). A lower signal to noise ratio in low expressed genes could explain this difference. However, another likely explanation pertains to the fact that co-transcriptional processing is not accounted for by the Zeisel model. While it has been documented that an RNA molecule is often processed while being synthesized (the "assembly-line model") [22], the Zeisel model considers synthesis and processing as two independent point events. This discrepancy is likely to be more relevant for short-lived (and thus low-expressed) transcripts, a sizeable fraction of which is expected to be

nascent at sequencing time. Those nascent transcripts may contribute to an intron to exon ratios higher than one when they are incompletely synthesized (for example if the last exon has not yet been produced). This hypothesis is corroborated by considering unspliced transcripts length, which putatively affects synthesis time and thus the probability of being nascent at sequencing time. Transcripts lying above the $r_l > 1$ line are indeed longer than those lying below this line (p – value $< 10^{-100}$, Wilcoxon test).

The transcripts incompatible with our model, amount to 25% of protein-coding genes with an exon TPM higher than 1, and are discarded from further analyses. The processing and degradation rates were computed either by solving (12) when $r_l > (2 - r_u)^{-1}$ or by optimizing (15) otherwise. For these cases that had two solutions (6% of the transcripts), we selected the one corresponding to rates most consistent with the other transcripts.

The resulting synthesis and processing rates are depicted in Fig. 3, right. Although processing rates span a smaller range of values, they are highly correlated (72%), which is not surprising as RNA processing occurs co-transcriptionally [22]. More remarkable is the correlation of synthesis and degradation rates, displayed in Fig. 4, left. At 65%, it is very similar to the 66% reported by [21] for the same cell type. This is also consistent with the emerging concept of a coupling between RNA transcription and decay [28]. Our data indicate that genes span a large range of dynamics, irrespective of their expression level. Indeed, genes with high synthesis and degradation rates can have the same steady-state expression level as genes with low synthesis and degradation rates. However, the former will reach this steady state faster than the latter. It thus makes sense to consider our RNA metabolic rates in the functional frame of reference indicated in Fig. 4, left. One axis corresponds to the steady state RNA abundance, given by the log-ratio of synthesis over degradation rates (or equivalently by the difference of log of the rates). The second axis correspond to the responsiveness of the gene, i.e. how fast it reaches steady state (computed by the sum of the log of the synthesis and degradation rates). It has been observed before that genes involved in more reactive and dynamic biological processes such as chromatin remodeling or transcription regulation tend to have a higher turnover than genes involved in more stable processes such as basic metabolism [9]. We checked that our data confirm this observation by looking at the Gene Ontology (GO,[29]) annotations of biological processes most associated by [9] with high and low turnover, namely "transcription" and "monosaccharide metabolism". Despite having similar steady-state abundances, transcripts of genes involved in transcription indeed have significantly faster dynamics and the ones involved in monosaccharide metabolism have significantly slower dynamics than the rest of the genes, as illustrated by the squares in Fig. 4, left and right. Other categories where our data confirms faster genes include chromatin modifications, cell cycle and transcription regulation.

We assessed the precision of our method by comparing the resulting degrada-

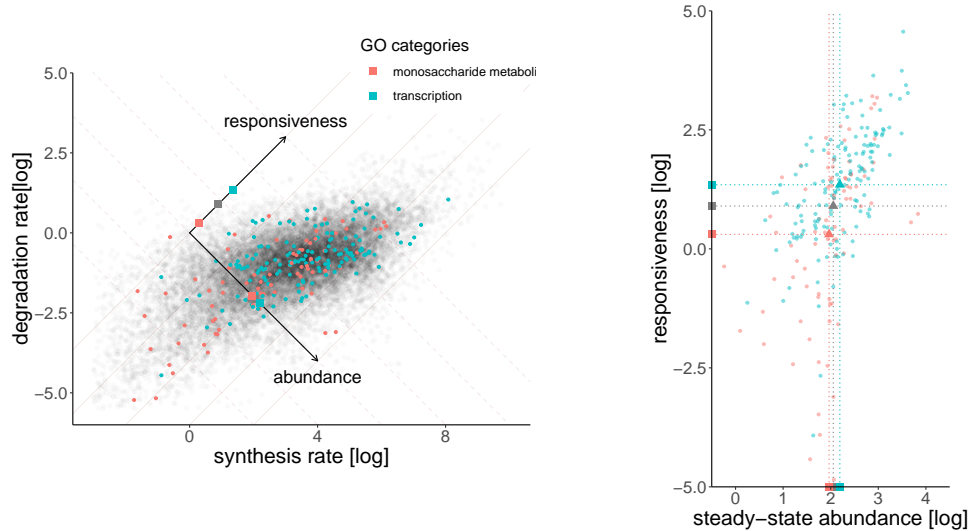


Figure 4: Left: Estimated RNA synthesis and degradation rates obtained from a single sample. These rates can also be considered in a different and maybe functionally more relevant frame of reference defined by the steady state abundance (first axis) and gene responsiveness (second and perpendicular axis), as illustrated by the background grid. Genes involved in fast adapting biological processes (such as transcription) tend to be more responsive than genes involved in stable functions (such as monosaccharide metabolism). The squares on the axes represent the projections of the mean rates for the respective categories (gray representing genes that belong to neither of the two categories) and indicate that mean transcript responsiveness (but not abundance) is strongly affected by the category. These two GO categories were selected for illustration because they were previously reported to be mostly enriched in high and low turn-over genes respectively [9]. Right: Same data as in left, but rotated and showing only colored dots, for visibility.

tion rates to those published for the same cell type by [21]. Those were obtained by using three replicates and seven time points and applying the SLAM-seq nucleotide-conversion method that, unlike metabolic labeling, does not require biochemical separation between the labeled and unlabeled RNA and is thus not affected by noise generated by the imperfect separation process (although that method has its own source of noise). From our data, we obtained gene degradation rates by taking, for each gene, the weighted average degradation rates of the corresponding transcripts. The weights were given by the mean exonic expression levels (unlabeled and labeled). We expect a lower precision for transcripts close to the $r_l = 1$ line, for which the labeling time was likely somewhat too short, so to assess the correlation, we weighted the transcripts by $1 - r_l$. Fig. 5, left, compares degradation rates obtained in our experiments with those reported by [21], keeping only genes with an average expression value higher than 100 TPM. We expect a higher precision for highly expressed genes, as this allows for a more precise estimates of the intron to exon ratios. This is indeed the case, and depending on the expression threshold and the sample, the correlation between our data and the previously published rates, we obtain a correlation ranging between 30% and 67% for a single sample estimate (see Fig. 5, left). As these experiments were performed in different labs using different methods, these numbers show that our rates obtained on a single sample and time point are meaningful. For comparison, [30] report correlations around 70% by using the *same* data, but changing only the method of analysis. Using three replicates, [25] report a 26% correlation using the INSPEcT package.

3.3 Comparison with INSPEcT

Since our method estimates metabolic rates from a single sample, we decided to compare its results to the "initial guess" provided by the INSPEcT method, to our knowledge the only other method that does not need multiple samples. Note however, that those rates are only the initial step of the INSPEcT method, and should not be confused with the global outcome of INSPEcT, which then aggregates multiple samples for the estimation. For concision, we will in the section refer to our method as SSRE (Single Sample Rate Estimation). The main differences between the two approaches is that INSPEcT assumes no degradation on labeled RNA and requires the estimation of a scaling factor accounting for the difference in RNA concentration between labeled and unlabeled samples, which SSRE avoids by considering the intron to exon ratio in each sample separately for each sample. Furthermore, INSPEcT requires the estimation the time derivative of the RNA abundances, which is avoided in SSRE by taking advantage of the analytical solution to the Zeisel model.

We used the INSPEcT package for R (more specifically the `newINSPEcT` function with parameter `preexisting=TRUE` and the `ratesFirstGuess` function) on the same data. The expression variance required by this function was estimated from the expression level from all three replicates using Loess regression. It took about 10 minutes to estimate rates for each replicate (about 50000 transcripts)

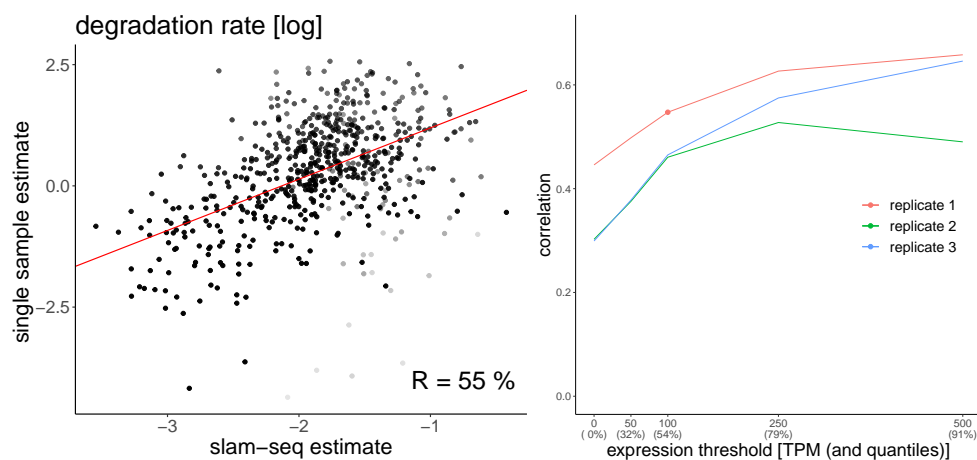


Figure 5: Left: degradation rates estimated from a single sample plotted against degradation rates published in [21] (obtained using slam-seq). The red line is obtained through weighted linear regression. The weights are set as $1 - r_l$ as indicated by the transparency of the dots. The (weighted) correlation of 55% indicates that the estimated rates are meaningful. Only genes with a mean exon TPM above 100 are taken into account. Right: Correlation between degradation rates obtained by [21] and the ones obtained our single-sample method as a function of expression level. Each line represents a biological replicate. The dot corresponds to the data shown on the left. As expected, the correlation is higher for highly expressed genes, as the intro to exon ratios can be more reliably estimated. In this experiment, replicate 1 correlates better than the two others, indicating that it is probably of better quality.

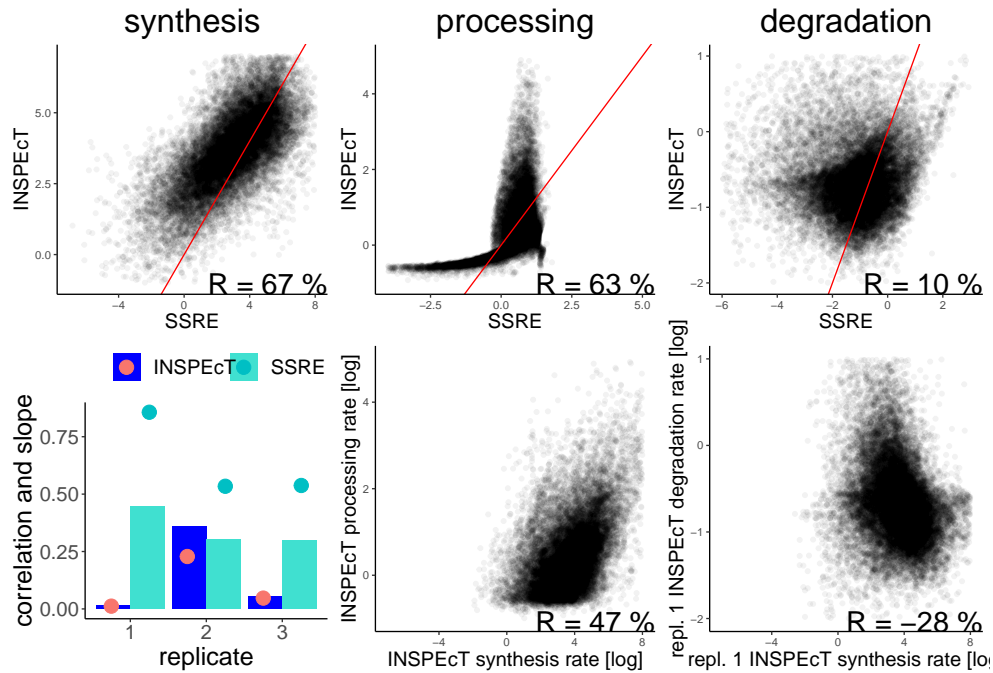


Figure 6: Comparison of our method (SSRE) with the INSPEcT "first guess" on the same data. Top row: direct comparison of rates obtained with our method and with the INSPEcT package on a single sample. Synthesis and processing rates are well correlated but not the degradation rate (Spearman correlation shown). The red bar indicates the diagonal. Bottom left: bars indicate the correlation of degradation rates with previously published data [21], as in Fig. 5. The INSPEcT method provides degradation rates with good correlation only for one of the three replicates (repl. 2), whereas it is the case for all three replicates using our method. The dots indicate the slope of the regression line in log-log space (as in Fig. 5, left). Slopes obtained from SSRE estimates are closer to one, which correspond to the ideal case of a linear relationship between the (non-log) rates. Bottom center and right: Rates obtained with INSPEcT also reproduce the positive correlation between synthesis and processing rates, but they produce a negative correlation between synthesis and degradation rates, unlike our method and previously published results.

using a single 2.8GHz core from a laptop computer, whereas our implementation of SSRE took about 30 seconds to complete on the same machine.

In addition to direct rate comparison, we decided to compare the methods using three criteria: (1) correlation with published rate, (2) rate distributions and (3) reproducibility across replicates. Figure 6 shows that the two methods provides synthesis and processing rates that are well correlated, while degradation rates are not. Moreover, the degradation rates obtained by INSPEcT correlate well to previously published rates only for one of the three samples. In contrast to degradation rates obtained with SSRE, the correlation with previously published rates does not improve when focusing on highly expressed genes, it even become negative for replicates 1 and 3 (data not shown). This suggests that for SSRE, degradation rate estimation is likely to improve with higher sequencing depth (and thus a more precise estimate of the intron to exon ratio). Finally, the rates computed using the INSPEcT method do not exhibit the previously documented correlation between synthesis and degradation rates [21]. This leads us to think that our degradation rates are closer to the real rates than the ones provided by the INSPEcT "first guess". This should not come as a surprised, as our method does not assume that labeled RNA does not degrade and estimates degradation and processing rates independently from the synthesis rate.

The synthesis and processing rates provided by the two methods are relatively well correlated, and INSPEcT provides rates that are more consistent across replicates (see Supplementary Fig D.8.) It is also interesting to note that SSRE tends to show an upper bound for the processing rate, while INSPEcT first shows a lower bound for that rate. It is difficult to speculate which (if any) is more likely true, but an upper bound would be consistent with biophysical constraints in a leaky co-transcriptional RNA splicing setting. Fig D.8 also shows that unlike INSPEcT, SSRE computes degradation rates that span a larger range of values than processing rates, a property also reported in [11] for a different system.

4 Discussion

In this paper, we presented a method to estimate synthesis, processing and degradation rates of RNA transcripts from a single 4sU labeled sample. We validated our method first *in silico* and then on real data obtained from mouse embryonic stem cells. Using our method we first replicated, on a different cell type, previous findings about the enrichment in high or low turn-over genes of specific cellular processes. Second, we showed that the rates obtained with our method correlate well (between 30% and 67%) with published rates obtained by applying SLAM-seq to the same cell types. Methods for such estimation have been published before, but they usually require a sufficient number of samples (around a dozen). We compared our method to the initial step of the INSPEcT method, which handles each sample separately, and obtain similar synthesis and processing rates, but different degradation rates. Our rates correlate more

consistently with previously published degradation rates obtained with nuclear conversion methods on the same system, and even more so for highly expressed transcripts. Rates obtain with our method also better reproduce previously observed statistical relationships between rates, although synthesis and processing rates are less consistent across replicates. Taken together these results suggest that our method provides more reliable degradation rates.

In contrast to other methods, our method explicitly uses the analytical solution to the Zeisel model of RNA dynamics. Moreover, our method is self-normalizing as it only uses the ratio of intron to exon expression levels. It is thus not affected by differences in sequencing depth of the various samples. This approach makes our method also faster than other methods as it boils down to numerically solving on a bounded domain either a univariate equation or a one-dimensional optimization for each transcript. Our method could thus be a suitable alternative to the initial step of the INSPEcT method especially when using a large number of samples as it is also about 20 times faster.

Similarly to the initial step of the INSPEcT method, a caveat of our method is that a sizable fraction of mostly lowly expressed transcripts (about 25 % in our case) are inconsistent with the model and their dynamics cannot be estimated. Together with the high correlation between synthesis and processing rate, it suggests that modeling transcription and processing as independent events is a simplification that could be reconsidered, as the coupling between the two has been documented [22]. However, this limitation of the Zeisel model is likely to also affect other methods using it [20, 31].

Another limitation of the method is that, unlike in [20], it does not consider the effect of leakage of unlabeled RNA in the labeled RNA pool because of unspecific capture. This leakage has the effect of reducing r_l towards the diagonal, and could potentially be estimated from the data as it is shared across all transcripts. Another improvement would be to embed this method in a probabilistic framework in order to quantify the estimate uncertainty (as in [30] for a simpler model) or to determine the optimal labeling time (as in [32]).

While using a single sample allows to reduce costs, this is not the only merit of this approach. In practice most experiments will have biological replicates, in which case our methods enables obtaining point estimates of α , β and γ for each of them. This in turn allows for estimating their variance, as well as assessing sample quality (e.g. if one of them systematically gives very different estimates for all genes). Moreover, because cell growth is likely to be limited during (short) labeling time, it is less likely to interfere in the estimation process than when using time course data, where it can have an effect [18]. In addition, when used in a time-course experiment, our method allows to investigate the evolution of these rates over time and assess whether these rates are stationary. Finally, the theoretical results obtained in this paper, could be used to improve other methods. For example, the method could be used to analyze SLAM-seq data which would reduce the number of samples but also provide estimate for the processing rate. Another possible application is single cell RNA velocity [31],

where the Zeisel model of RNA dynamics is also used, but splicing rates γ are set to be equal for all transcripts. While it has been documented (and is consistent with our data) that splicing rates are more homogeneous than degradation rates [11], this is potentially an approximation that could be improved with our framework to increase the accuracy of the method, for example by considering the strong correlation between the synthesis and processing rates. Finally, our method could also be used in conjunction with the recent developments in single cell metabolic labeling experiments [33, 34].

Acknowledgements

This work was funded by the Swiss National Science Foundation through grant no. FN 310030_152724/1 to S.B and PP00P3_150667 and the NCCR in RNA & Disease to A.C.M.

Competing interests

The authors declare that they have no conflicts of interest.

Availability of data and code

An R package implementing our method is available on github, together with the code used to generate the figures as well as the gene expression data used: <https://github.com/BergmannLab/SingleSampleRNAdynamics>

The raw data files data are available on the Gene Expression Omnibus accession number GEO:GSE150286 (main replicate) and GEO: GSE143277 (second and third replicates of Fig. 5, right and Supplementary Fig. D.8).

Author's contribution

M.H. developed and implemented the method, analyzed the expression quantification data, interpreted the results, figured out the proof, generated the figures and wrote the manuscript, A.B. performed the experiments and interpreted the results, A.C.M. initiated the project, designed the study and generated the expression quantification data, S.B. reviewed the math, interpreted the results and supervised the process. All authors contributed to the manuscript.

References

- [1] Petit, F., Sears, K.E., Ahituv, N.: Limb development: a paradigm of gene regulation. *Nature Reviews Genetics* **18**(4), 245 (2017)

- [2] Lee, T.I., Young, R.A.: Transcriptional regulation and its misregulation in disease. *Cell* **152**(6), 1237–1251 (2013)
- [3] Mack, K.L., Ballinger, M.A., Phifer-Rixey, M., Nachman, M.W.: Gene regulation underlies environmental adaptation in house mice. *Genome research* **28**(11), 1636–1645 (2018)
- [4] Elkon, R., Zlotorynski, E., Zeller, K.I., Agami, R.: Major role for mrna stability in shaping the kinetics of gene induction. *BMC genomics* **11**(1), 259 (2010)
- [5] Alpert, T., Herzelt, L., Neugebauer, K.M.: Perfect timing: splicing and transcription rates in living cells. *WIREs: RNA* **8**(2), 1401 (2017)
- [6] Friedel, C.C., Dölken, L.: Metabolic tagging and purification of nascent rna: implications for transcriptomics. *Molecular BioSystems* **5**(11), 1271–1278 (2009)
- [7] Dölken, L., Ruzsics, Z., Rädle, B., Friedel, C.C., Zimmer, R., Mages, J., Hoffmann, R., Dickinson, P., Forster, T., Ghazal, P., *et al.*: High-resolution gene expression profiling for simultaneous kinetic parameter analysis of RNA synthesis and decay. *Rna* **14**(9), 1959–1972 (2008)
- [8] Rabani, M., Levin, J.Z., Fan, L., Adiconis, X., Raychowdhury, R., Garber, M., Gnirke, A., Nusbaum, C., Hacohen, N., Friedman, N., *et al.*: Metabolic labeling of rna uncovers principles of rna production and degradation dynamics in mammalian cells. *Nature biotechnology* **29**(5), 436 (2011)
- [9] Schwanhäusser, B., Busse, D., Li, N., Dittmar, G., Schuchhardt, J., Wolf, J., Chen, W., Selbach, M.: Global quantification of mammalian gene expression control. *Nature* **473**(7347), 337 (2011)
- [10] Windhager, L., Bonfert, T., Burger, K., Ruzsics, Z., Krebs, S., Kaufmann, S., Malterer, G., L’Hernault, A., Schilhabel, M., Schreiber, S., *et al.*: Ultrashort and progressive 4su-tagging reveals key characteristics of RNA processing at nucleotide resolution. *Genome research* **22**(10), 2031–2042 (2012)
- [11] Rabani, M., Raychowdhury, R., Jovanovic, M., Rooney, M., Stumpo, D.J., Pauli, A., Hacohen, N., Schier, A.F., Blackshear, P.J., Friedman, N., *et al.*: High-resolution sequencing and modeling identifies distinct dynamic RNA regulatory strategies. *Cell* **159**(7), 1698–1710 (2014)
- [12] Barrass, J.D., Reid, J.E., Huang, Y., Hector, R.D., Sanguinetti, G., Beggs, J.D., Granneman, S.: Transcriptome-wide RNA processing kinetics revealed using extremely short 4tu labeling. *Genome biology* **16**(1), 282 (2015)
- [13] Furlan, M., de Pretis, S., Pelizzola, M.: Dynamics of transcriptional and post-transcriptional regulation. *Briefings in Bioinformatics* (2020)

- [14] Neymotin, B., Athanasiadou, R., Gresham, D.: Determination of in vivo rna kinetics using rate-seq. *Rna* **20**(10), 1645–1652 (2014)
- [15] Uvarovskii, A., Dieterich, C.: pulseR: Versatile computational analysis of rna turnover from metabolic labeling experiments. *Bioinformatics* **33**(20), 3305–3307 (2017)
- [16] Lugowski, A., Nicholson, B., Rissland, O.S.: Determining mRNA half-lives on a transcriptome-wide scale. *Methods* **137**, 90–98 (2018)
- [17] Russo, J., Heck, A.M., Wilusz, J., Wilusz, C.J.: Metabolic labeling and recovery of nascent RNA to accurately quantify mRNA stability. *Methods* **120**, 39–48 (2017)
- [18] Lugowski, A., Nicholson, B., Rissland, O.S.: DRUID: A pipeline for transcriptome-wide measurements of mRNA stability. *RNA* **24**(5), 623–632 (2018)
- [19] De Pretis, S., Kress, T., Morelli, M.J., Mellon i, G.E., Riva, L., Amati, B., Pelizzola, M.: INSPEcT: a computational tool to infer mrna synthesis, processing and d egradation dynamics from rna-and 4su-seq time course experiments. *Bioinformatics* **31**(17), 2829–2835 (2015)
- [20] Furlan, M., Galeota, E., Del Gaudio, N., Dassi, E., Caselle, M., de Pretis, S., Pelizzola, M.: Genome-wide dynamics of rna synthesis, processing, and degradation without rna metabolic labeling. *Genome Research* **30**(10), 1492–1507 (2020)
- [21] Herzog, V.A., Reichholf, B., Neumann, T., Rescheneder, P., Bhat, P., Burkard, T.R., Wlotzka, W., von Haeseler, A., Zuber, J., Ameres, S.L.: Thiol-linked alkylation of RNA to assess expression dynamics. *Nature methods* **14**(12), 1198 (2017)
- [22] Herzel, L., Ottoz, D.S., Alpert, T., Neugebauer, K.M.: Splicing and transcription touch base: co-transcriptional spliceosome assembly and function. *Nature reviews Molecular cell biology* **18**(10), 637 (2017)
- [23] Zeisel, A., Köstler, W.J., Molotski, N., Tsai, J.M., Krauthgamer, R., Jacob-Hirsch, J., Rechavi, G., Soen, Y., Jung, S., Yarden, Y., et al.: Coupled pre-mRNA and mRNA dynamics unveil operational strategies underlying transcriptional responses to stimuli. *Molecular systems biology* **7**(1) (2011)
- [24] Biasini, A., Abdulkarim, B., de Pretis, S., Tan, J.Y., Arora, R., Wischniewski, H., Dreos, R., Pelizzola, M., Ciaudo, C., Marques, A.C.: Translation is required for mirna-dependent decay of endogenous transcripts. *The EMBO journal* **40**(3), 104569 (2021)
- [25] Biasini, A., Marques, A.C.: A protocol for transcriptome-wide inference of RNA metabolic rates in mouse embryonic stem cells. *Frontiers in Cell and Developmental Biology* **8**, 97 (2020)

- [26] Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., Gingeras, T.R.: Star: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**(1), 15–21 (2013)
- [27] Li, B., Dewey, C.N.: RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC bioinformatics* **12**(1), 323 (2011)
- [28] Hartenian, E., Glaunsinger, B.A.: Feedback to the central dogma: cytoplasmic mrna decay and transcription are interdependent processes. *Critical reviews in biochemistry and molecular biology* **54**(4), 385–398 (2019)
- [29] Consortium, G.O.: The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic acids research* **47**(D1), 330–338 (2019)
- [30] Jürges, C., Dölken, L., Erhard, F.: Dissecting newly transcribed and old RNA using GRAND-SLAM. *Bioinformatics* **34**(13), 218–226 (2018)
- [31] La Manno, G., Soldatov, R., Zeisel, A., Braun, E., Hochgerner, H., Petukhov, V., Lidschreiber, K., Kastrioti, M.E., Lönnerberg, P., Furlan, A., *et al.*: RNA velocity of single cells. *Nature* **560**(7719), 494 (2018)
- [32] Uvarovskii, A., Naarmann-de Vries, I.S., Dieterich, C.: On the optimal design of metabolic RNA labeling experiments. *PLoS computational biology* **15**(8), 1007252 (2019)
- [33] Battich, N., Beumer, J., de Barbanson, B., Krenning, L., Baron, C.S., Tanenbaum, M.E., Clevers, H., van Oudenaarden, A.: Sequencing metabolically labeled transcripts in single cells reveals mrna turnover strategies. *Science* **367**(6482), 1151–1156 (2020)
- [34] Cao, J., Zhou, W., Steemers, F., Trapnell, C., Shendure, J.: Sci-fate characterizes the dynamics of gene expression in single cells. *Nature biotechnology* **38**(8), 980–988 (2020)

A Derivation of the model solution

This is a first order linear ordinary differential equation in $p(t)$ and $m(t)$ that can be expressed in matrix form as

$$\begin{pmatrix} \dot{p} \\ \dot{m} \end{pmatrix} = \begin{pmatrix} -\beta & 0 \\ \beta & -\gamma \end{pmatrix} \begin{pmatrix} p \\ m \end{pmatrix} + \begin{pmatrix} \alpha \\ 0 \end{pmatrix} \quad (17)$$

The solution to this equation is given by

$$\begin{pmatrix} p \\ m \end{pmatrix} = k_1 \mathbf{v} \exp(\lambda_1 t) + k_2 \mathbf{w} \exp(\lambda_2 t) + \begin{pmatrix} \frac{\alpha}{\beta} \\ \frac{\alpha}{\gamma} \end{pmatrix}, \quad (18)$$

where k_1 and k_2 are scalar constants determined by the boundary conditions, λ_1 , λ_2 are eigenvalues of the matrix in (17) and \mathbf{v} , \mathbf{w} are the corresponding eigenvectors.

The eigenvalues are given by $\lambda_1 = -\beta$ and $\lambda_2 = -\gamma$. The first eigenvector \mathbf{v} is obtained by solving

$$\begin{cases} -\beta \mathbf{v}_1 = -\beta \mathbf{v}_1 \\ \beta \mathbf{v}_1 - \gamma \mathbf{v}_2 = -\beta \mathbf{v}_2 \end{cases} \Rightarrow \mathbf{v}_1 = \frac{\gamma - \beta}{\beta} \mathbf{v}_2 \Rightarrow \mathbf{v} \propto \begin{pmatrix} \gamma - \beta \\ \beta \end{pmatrix} \quad (19)$$

Similarly the second eigenvector is obtained by solving

$$\begin{cases} -\beta \mathbf{w}_1 = -\gamma \mathbf{w}_1 \\ \beta \mathbf{w}_1 - \gamma \mathbf{w}_2 = -\gamma \mathbf{w}_2 \end{cases} \Rightarrow \mathbf{w}_1 = 0 \Rightarrow \mathbf{w} \propto \begin{pmatrix} 0 \\ 1 \end{pmatrix}.$$

The solution to (17) is thus given by

$$\begin{pmatrix} p \\ m \end{pmatrix} = k_1 \begin{pmatrix} \gamma - \beta \\ \beta \end{pmatrix} \exp(-\beta t) + k_2 \begin{pmatrix} 0 \\ 1 \end{pmatrix} \exp(-\gamma t) + \begin{pmatrix} \frac{\alpha}{\beta} \\ \frac{\alpha}{\gamma} \end{pmatrix}.$$

Expressed by its component this is equivalent to

$$p(t) = k_1(\gamma - \beta) \exp(-\beta t) + \frac{\alpha}{\beta} \quad (20)$$

$$m(t) = k_1 \beta \exp(-\beta t) + k_2 \exp(-\gamma t) + \frac{\alpha}{\gamma} \quad (21)$$

We now turn to the boundary conditions to determine k_1 and k_2 . The boundary conditions are different for the unlabeled and the labeled RNA.

Unlabeled RNA

Like in [19], we assume the system to be in steady-state prior to labeling. The steady-state is given by solving (17) with $\dot{p} = \dot{m} = 0$.

$$\begin{cases} 0 = -\beta p + \alpha \\ 0 = \beta p - \gamma m \end{cases} \Rightarrow \begin{cases} p = \frac{\alpha}{\beta} \\ 0 = \beta \frac{\alpha}{\beta} - \gamma m \end{cases} \Rightarrow \begin{cases} p = \frac{\alpha}{\beta} \\ m = \frac{\alpha}{\gamma} \end{cases} \quad (22)$$

During labeling time, we assume that no unlabeled RNA is synthesized such that $\alpha = 0$. Assuming that we start labeling at time $t = 0$, we thus have

$$p_u(0) = \frac{\alpha}{\beta} \Rightarrow k_1(\gamma - \beta) = \frac{\alpha}{\beta} \Rightarrow k_1 = \frac{\alpha}{\beta(\gamma - \beta)} \quad (23)$$

Moreover we have

$$m_u(0) = \frac{\alpha}{\gamma} \Rightarrow \frac{\alpha}{\gamma - \beta} + k_2 = \frac{\alpha}{\gamma} \Rightarrow k_2 = \frac{\alpha}{\gamma} - \frac{\alpha}{\gamma - \beta} = \frac{-\beta\alpha}{\gamma(\gamma - \beta)}$$

This leads us to the solution for the unlabeled RNA

$$p_u(t) = \frac{\alpha}{\beta} \exp(-\beta t) \quad (24)$$

$$m_u(t) = \frac{\alpha}{\gamma - \beta} \exp(-\beta t) - \frac{\beta\alpha}{\gamma(\gamma - \beta)} \exp(-\gamma t), \quad (25)$$

where the u label indicates that this corresponds to the unlabeled RNA pool.

Labeled RNA

The solution for the labeled RNA could be obtained the same way as for the unlabeled RNA, but setting $\alpha \neq 0$ and $p_l(0) = m_l(0) = 0$. However, it is simpler to notice that the total RNA (labeled and non-labeled) stay at steady-state during the labeling such that we have the following solution for labeled RNA.

$$p_l(t) = \frac{\alpha}{\beta} - p_u(t) = \frac{\alpha}{\beta} (1 - \exp(-\beta t))$$

$$m_l(t) = \frac{\alpha}{\gamma} - m_u(t) = \frac{\alpha}{\gamma} \left(1 + \frac{\beta}{(\gamma - \beta)} \exp(-\gamma t)\right) - \frac{\alpha}{\gamma - \beta} \exp(-\beta t)$$

where the l label indicates that this corresponds to the labeled RNA pool.

B Equation simplification

In this appendix we show how the the system given by equations (10, 11) can be simplified to Eq. (12) to infer the ratio k between the processing and degradation rate. Starting from

$$a = \frac{(1 - k)E_{k\gamma}}{E_{k\gamma} - k^2 E_\gamma} \quad (26)$$

$$b = \frac{(1 - k)(1 - E_{k\gamma})}{(1 - E_{k\gamma}) - k^2(1 - E_\gamma)}, \quad (27)$$

we have

$$a(E_{k\gamma} - k^2 E_\gamma) = (1 - k)E_{k\gamma} \quad (28)$$

$$b((1 - E_{k\gamma}) - k^2(1 - E_\gamma)) = (1 - k)(1 - E_{k\gamma}). \quad (29)$$

Summing (28) and (29) yields

$$E_{k\gamma}(a-b) + k^2 E_\gamma(b-a) + b(1-k^2) = 1-k \quad (30)$$

$$\begin{aligned} \Leftrightarrow E_{k\gamma} - k^2 E_\gamma &= \frac{(1-k) - b(1-k^2)}{a-b} \\ &= \frac{(1-k)(1-b(1+k))}{a-b}. \end{aligned} \quad (31)$$

Dividing (28) by (29) and inserting (31) results in

$$\begin{aligned} \frac{E_{k\gamma}}{1-E_{k\gamma}} &= \frac{a}{b} \frac{E_{k\gamma} - k^2 E_\gamma}{(1-E_{k\gamma}) - k^2(1-E_\gamma)} \\ &= \frac{a}{b} \frac{E_{k\gamma} - k^2 E_\gamma}{(1-k^2) - (E_{k\gamma} - k^2 E_\gamma)} \\ &= \frac{a}{b} \frac{(1-k)(1-b(1+k))}{(1-k^2)(a-b) - (1-k)(1-b(1+k))} \\ &= \frac{a}{b} \frac{1-b(1+k)}{(1+k)(a-b) - 1 + b(1+k)} \\ &= \frac{a}{b} \frac{1-b(1+k)}{(1+k)a-1} = -\frac{a-ab(1+k)}{b-ab(1+k)} \end{aligned} \quad (32)$$

It follows that

$$E_{k\gamma}(b-ab(1+k)) = (E_{k\gamma}-1)(a-ab(1+k)) \quad (33)$$

$$\Leftrightarrow (b-a)E_{k\gamma} = ab(1+k) - a, \quad (34)$$

and thus

$$\exp(-k\gamma T) = E_{k\gamma} = \frac{kab + ab - a}{b-a} = \frac{a(bk + b - 1)}{b-a} \quad (35)$$

$$\Leftrightarrow k\gamma T = \log\left(\frac{b-a}{a(bk + b - 1)}\right) \quad (36)$$

Moreover, from (10), we have that

$$\begin{aligned} a &= \frac{1-k}{1-k^2 \exp((k-1)\gamma T)} \Leftrightarrow \exp((k-1)\gamma T) = \frac{(1-\frac{1-k}{a})}{k^2} \\ \Leftrightarrow (k-1)\gamma T &= \log\left(\frac{k+a-1}{k^2 a}\right) \end{aligned} \quad (37)$$

Multiplying (36) by $\frac{k-1}{k}$ and subtracting (37) results in

$$0 = \frac{k}{k-1} \log\left(\frac{k+a-1}{k^2 a}\right) - \log\left(\frac{(b-a)}{a(bk + b - 1)}\right). \quad (38)$$

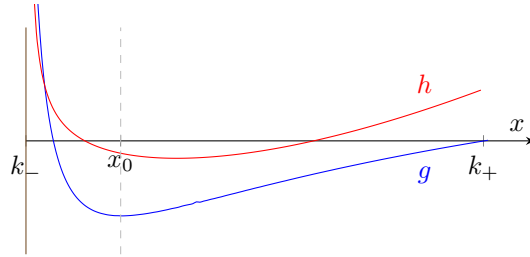


Figure C.7: Sketch of the proof that $g(x)$ has a single zero in $\mathcal{D} = (k_-, k_+)$. We first show that $\lim_{x \rightarrow k_-} g(x) = \infty$, that $\lim_{x \rightarrow k_+} g(x) = 0$ and $\lim_{x \rightarrow k_+} g'(x) > 0$, so that g must cross the x-axis on \mathcal{D} . To show that it only does it once, we consider a function $h(x)$ that has the same sign as $g(x)$ when $g'(x) = 0$. We show that h is convex on \mathcal{D} and thus g cannot have a negative extrema, followed by a positive extrema, followed by a negative extrema. Hence it cannot have more than one zero on \mathcal{D} .

This equation also provides upper and lower bounds for k as both $\frac{k+a-1}{a}$ and $bk + b - 1$ must be strictly positive for their logarithm to be defined and

$$0 < \exp(-\beta T) = E_{k\gamma} = \frac{kab + ab - a}{b - a} < 1 \quad \forall \beta T > 0 \quad (39)$$

for (35) to hold. Developing these three conditions results in the following domain of definition \mathcal{D} for k :

$$\max\left(\frac{1}{b} - 1, 1 - a\right) < k < \frac{1}{a} - 1, \quad (40)$$

where $0 < a < b < 1$.

C Proof of unicity of solution

In this appendix, we prove that (12) has a single solution for $b > (2 - a)^{-1}$. We first note that $b > (2 - a)^{-1} \Leftrightarrow \frac{1-b}{b} < 1 - a$, so the lower bound for k is $k_- = 1 - a$. We then define the right-hand side of (12) as

$$g(x) = \frac{x}{x-1} \log\left(\frac{1}{x^2} \left(1 - \frac{1-x}{a}\right)\right) - \log\left(\frac{(b-a)}{a(b(x+1)-1)}\right) \quad (41)$$

$$\begin{aligned} &= \frac{x}{x-1} \left(-2 \log(x) + \log(a+x-1) - \log(a) \right) \\ &\quad - \log(b-a) + \log(a) + \log(b(x+1)-1) \end{aligned} \quad (42)$$

We then observe that on that lower bound $\lim_{x \rightarrow 1-a} g(x) = +\infty$ because $a+x-1$ tends to zero and $x-1$ is negative.

On the other hand, for the upper bound $k_+ = \frac{1-a}{a}$, we have

$$a + x - 1 = \frac{a^2 + 1 - a - a}{a} = \frac{(1-a)^2}{a}$$

$$b(x+1) - 1 = \frac{b}{a} - 1 = \frac{b-a}{a}$$

we can deduce that the upper bound of x is a zero of g :

$$g\left(\frac{1-a}{a}\right) = \frac{1-a}{1-2a} \left(-2\log(1-a) + 2\log(a) + 2\log(1-a) - 2\log(a) \right) - \log(b-a) + \log(a) + \log(b-a) - \log(a) = 0. \quad (43)$$

Moreover, the derivative of g is given by:

$$g'(x) = -\frac{1}{(x-1)^2} \left(-2\log(x) + \log(a+x-1) - \log(a) \right) + \frac{x}{x-1} \left(-\frac{2}{x} + \frac{1}{a+x-1} \right) + \frac{b}{b(x+1)-1}. \quad (44)$$

Then

$$\begin{aligned} g'\left(\frac{1-a}{a}\right) &= \frac{-a^2}{(1-2a)^2} \left(-2\log(1-a) + 2\log(a) + 2\log(1-a) \right. \\ &\quad \left. - 2\log(a) \right) + \frac{1-a}{1-2a} \left(\frac{-2a}{1-a} + \frac{a}{(1-a)^2} \right) + \frac{ba}{b-a} \\ &= \frac{1}{1-2a} \frac{2a-1}{1-a} + \frac{ba}{b-a} \\ &= -\frac{a}{1-a} + \frac{ba}{b-a} \\ &= \frac{a^2(1-b)}{(b-a)(1-a)} > 0 \quad \text{if } 0 < a < b < 1. \end{aligned} \quad (45)$$

Since $g(k_+)$ reaches zero from below, while $g(k_-) > 0$, we can infer that $g(x)$ has a zero between k_- and k_+ as illustrated on Fig C.7.

To show that this zero is unique, we look at the sign of $g'(x)$.

We can rewrite $g'(x)$ as

$$g'(x) = \frac{1}{(x-1)^2} (A(x) - B(x)) \quad (46)$$

where

$$A(x) = \frac{(1-x)(x(ab+b-1) - (3b-2)(1-a))}{(x+a-1)(bx+b-1)}$$

$$B(x) = \log\left(\frac{a+x-1}{ax^2}\right)$$

Let x_0 be a zero of g' , i.e., the position of a local extrema of g . We have

$$\begin{aligned} g(x_0) &= \frac{x_0}{x_0 - 1} \log\left(\frac{x_0 + a - 1}{ax_0^2}\right) - \log\left(\frac{(b-a)}{a(b(x_0+1)-1)}\right) \\ &= \frac{-x_0(x_0(ab+b-1) - (3b-2)(1-a))}{(x_0+a-1)(bx_0+b-1)} \\ &\quad - \log\left(\frac{(b-a)}{a(b(x_0+1)-1)}\right) \end{aligned} \quad (47)$$

The second equality holds because $g'(x_0) = 0$ by definition of x_0 . By multiplying (47) by (bx_0+b-1) , which is positive, we can then define a new function $h(x)$ whose sign is the same as the sign of $g(x)$ for $x = x_0$ (see Fig C.7 for an illustration).

$$h(x) = C(x) - D(x), \quad (48)$$

where

$$\begin{aligned} C(x) &= \frac{-x(x(ab+b-1) - (3b-2)(1-a))}{(x+a-1)} \\ D(x) &= (bx+b-1) \log\left(\frac{(b-a)}{a(bx+b-1)}\right) \end{aligned} \quad (49)$$

We can now compute the second derivatives of $C(x)$ and $D(x)$.

$$\begin{aligned} C'(x) &= -\frac{(2x(ab+b-1) - (3b-2)(1-a))(x+a-1)}{(x+a-1)^2} \\ &\quad + \frac{x^2(ab+b-1) - x(3b-2)(1-a)}{(x+a-1)^2} \\ &= -\frac{x^2(ab+b-1) - 2x(1-a)(ab+b-1) + (3b-2)(1-a)^2}{(x+a-1)^2} \\ C''(x) &= -\frac{(2x(ab+b-1) - 2(1-a)(ab+b-1))(x+a-1)}{(x+a-1)^3} \\ &\quad + \frac{2(x^2(ab+b-1) - 2x(1-a)(ab+b-1) + (3b-2)(1-a)^2)}{(x+a-1)^3} \\ &= -\frac{2(1-a)^2(ab-2b+1)}{(x+a-1)^3} > 0 \quad \forall b > \frac{1}{2-a} \\ D'(x) &= b \log\left(\frac{b-a}{a}\right) - b(\log(bx+b-1) + 1) \\ D''(x) &= -\frac{b^2}{bx+b-1} < 0 \quad \forall x > k_- \end{aligned}$$

Hence

$$h''(x) = C''(x) - D''(x) < 0 \quad \forall x \in \mathcal{D}, \forall b > \frac{1}{2-a} \quad (50)$$

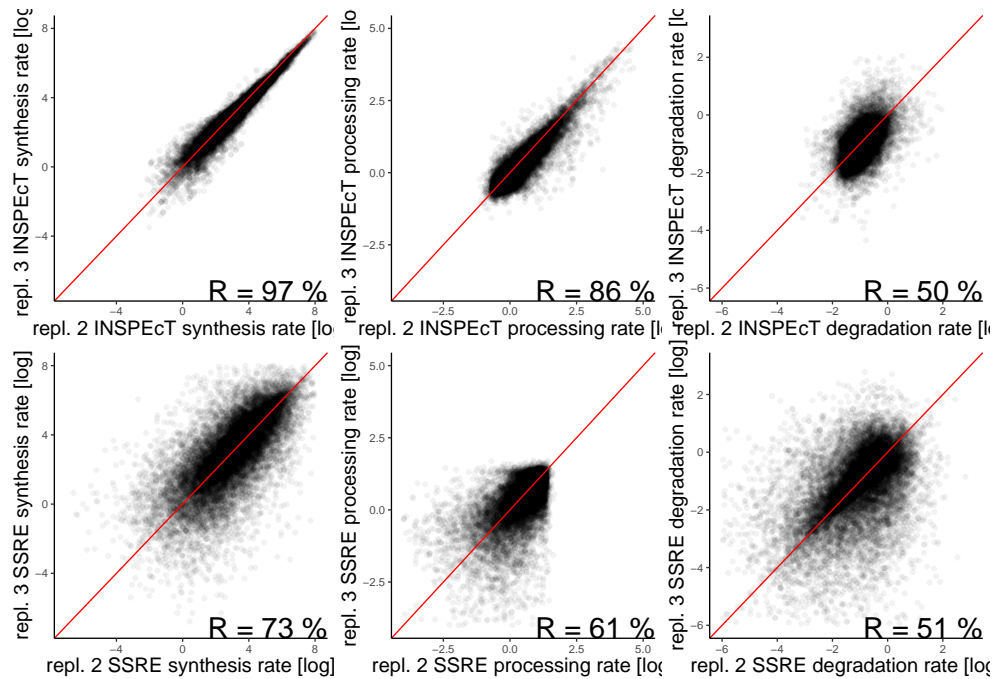


Figure D.8: Correlation of rates across two replicates. Top row: Rates obtained with INSPEcT. Bottom row: rates obtained with our method. The INSPEcT method provides rates that are more consistent across replicates for synthesis and processing rates, but not for degradation rates. Spearman correlation is indicated and the red line shows the diagonal.

This means that h is convex, so there cannot be three points $x_1 < x_2 < x_3$ such that $0 > h(x_1) < h(x_2) > 0 > h(x_3)$. Hence the same can be said of three zeros of g' , so $g(x)$ cannot have more than one zero. \square

D Supplementary figure