

Supporting Information

Aviran et al. 10.1073/pnas.1106541108

SI Text

Convexity and Simplification of the Optimization Problems. Before we start, we note that we replace n' with n throughout this section for notational convenience. To simplify the presentation, we also treat the *third* optimization subproblem first.

In the third optimization problem, we fix Γ and c at their most recent estimates, Γ^* and c^* , respectively, and seek a probability distribution Θ that maximizes $\log \mathcal{L}(\Theta, \Gamma^*, c^*)$. We call this problem **P3** and establish its strong convexity in the following lemma.

Lemma 1. *P3 is a feasible convex optimization problem.*

Proof: It is easy to see that maximizing $\log \mathcal{L}(\Theta, \Gamma^*, c^*)$ over Θ is equivalent to maximizing

$$\mathcal{F}(\Theta) = \sum_{k=1}^n X_k \left[c^* \left(\sum_{l=k+1}^n \theta_l - 1 \right) + \log(e^{c^* \theta_k} - (1 - \gamma_k^*)) \right]. \quad [\text{S1}]$$

We now show that $\mathcal{F}(\Theta)$ is strictly concave in Θ on its feasible convex domain $\mathcal{D}_{\mathcal{F}} = \{(\theta_1, \dots, \theta_n) | \theta_k > \frac{1}{c^*} \log(1 - \gamma_k^*) \forall 1 \leq k \leq n\}$. This can be seen from Eq. S1 by observing that $\mathcal{F}(\Theta)$ is a sum of a linear function of Θ and of log functions, where each logarithm depends only on one θ_k . For convenience, we denote the latter functions by $f_k(\Theta) = f_k(\theta_k) = \log(e^{c^* \theta_k} - (1 - \gamma_k^*))$ and infer their strict concavity from the sign of $f_k''(\theta_k) = \frac{-(1 - \gamma_k^*)(c^*)^2 e^{c^* \theta_k}}{f_k(\theta_k)^2} < 0$. We also wish to stress that $\mathcal{F}(\Theta)$'s domain includes negative θ_k 's because $\log(1 - \gamma_k^*) \leq \log(1)$, with a strict inequality for at least one k . Finally, the convexity and feasibility of **P3** follow directly from the linearity of the imposed equality and inequality constraints (1).

Next, we provide the complete details of the proof of Theorem 3 in the main text.

Proof of Theorem 3: For a feasible convex optimization problem, the Karush–Kuhn–Tucker (KKT) constraints provide necessary and sufficient conditions on its unique solution (1). From Lemma 1's proof, it follows that the KKT conditions for **P3** take the form

$$\begin{aligned} \theta_k^* \geq 0, \quad \lambda_k^* \geq 0, \quad \lambda_k^* \theta_k^* = 0, \quad k = 1, \dots, n, \\ \sum_{k=1}^n \theta_k^* = 1, \end{aligned} \quad [\text{S2}]$$

$$\frac{\partial \mathcal{F}(\Theta^*)}{\partial \theta_k} + \lambda_k^* - \nu^* = 0, \quad k = 1, \dots, n \quad [\text{S3}]$$

where Θ^* is the optimal solution, the λ_k^* 's and ν^* are the KKT multipliers that solve Eqs. S2 and S3, and $\mathcal{F}(\Theta)$ is defined in Lemma 1's proof. Eq. S3 can then be explicitly written as

$$c^* \left[\sum_{i=1}^{k-1} X_i + \frac{X_k}{1 - (1 - \gamma_k^*) e^{-c^* \theta_k^*}} \right] + \lambda_k^* = \nu^*, \quad 1 \leq k \leq n. \quad [\text{S4}]$$

Clearly, if the solution of the original, less-constrained problem is nonnegative, then it is also the sought solution and satisfies the KKT conditions with $\lambda_k^* = 0$ for all k .

Now, assume without loss of generality (*w.l.o.g.*) that $X_k > 0$. Consider the case where $\theta_k^* > 0$, implying that $\lambda_k^* = 0$ and that the following equation must hold:

$$c^* \sum_{i=1}^{k-1} X_i + X_k f_k'(\theta_k^*) = \nu^*, \quad [\text{S5}]$$

where we use the previously introduced notation $f_k(\theta) = \log(e^{c^* \theta} - (1 - \gamma_k^*))$ and where $f_k'(\theta)$ stands for $f_k(\theta)$'s derivative. The left-hand side of Eq. S5 consists of a nonnegative constant and of a strictly monotonously decreasing function of θ_k^* , as follows from Lemma 1's proof and as illustrated in Fig. S1A. For $\theta_k^* > 0$, $f_k'(\theta_k^*)$ is bounded by $[f_k'(\infty), f_k'(0)] = [c^*, \frac{c^*}{\gamma_k^*}]$, and, therefore, given $\nu^* \in (c^* \sum_{i=1}^k X_i, c^* (\sum_{i=1}^{k-1} X_i + \frac{X_k}{\gamma_k^*}))$, one can solve Eq. S5 for a positive θ_k^* . The solution then satisfies

$$e^{c^* \theta_k^*} = (1 - \gamma_k^*) \frac{\nu^* - c^* \sum_{i=1}^{k-1} X_i}{\nu^* - c^* \sum_{i=1}^k X_i} > 1. \quad [\text{S6}]$$

However, if ν^* exceeds the value of the right boundary point of the interval $[c^*, \frac{c^*}{\gamma_k^*}]$, then $\theta_k^* > 0$ is impossible, and thus we set $\theta_k^* = 0$ and $\lambda_k^* = \nu^* - c^* (\sum_{i=1}^{k-1} X_i + \frac{X_k}{\gamma_k^*}) \geq 0$. Note that if $\gamma_k^* = 0$, $f_k'(\theta_k^*)$ is unbounded over \mathbb{R}^+ and so any $\nu^* > c^* \sum_{i=1}^k X_i$ yields a positive solution. Also note that the explicit expression for $\hat{\Theta}$ that is computed by Algorithm 1 is obtained when plugging $\nu^* = \hat{c} \sum_{i=1}^n X_i$ into Eq. S6, with $c^* = \hat{c}$ and $\gamma_k^* = \hat{\gamma}_k$.

For notational convenience, we introduce the constants $\alpha_k = c^* (\sum_{i=1}^{k-1} X_i + \frac{X_k}{\gamma_k^*})$, $1 \leq k \leq n$, and obtain the following relations between $e^{c^* \theta_k^*}$ and ν^* :

$$e^{c^* \theta_k^*} = \begin{cases} 1 & \text{if } \nu^* \geq \alpha_k \\ (1 - \gamma_k^*) \frac{\nu^* - c^* \sum_{i=1}^{k-1} X_i}{\nu^* - c^* \sum_{i=1}^k X_i} & \text{if } c^* \sum_{i=1}^k X_i < \nu^* < \alpha_k \\ \text{undefined} & \text{if } \nu^* \leq c^* \sum_{i=1}^k X_i \end{cases} \quad [\text{S7}]$$

Note that when $\nu^* \leq c^* \sum_{i=1}^k X_i$, there is no feasible solution to the KKT conditions because the θ_k^* that solves Eq. S5 lies outside of $\mathcal{F}(\Theta)$'s domain. As will become clear later, this regime is irrelevant to the problem's solution. Hence, we restrict attention to $\nu^* > c^* \sum_{i=1}^k X_i$ and simplify Eq. S7 to

$$e^{c^* \theta_k^*} = \max \left\{ 1, (1 - \gamma_k^*) \frac{\nu^* - c^* \sum_{i=1}^{k-1} X_i}{\nu^* - c^* \sum_{i=1}^k X_i} \right\}. \quad [\text{S8}]$$

Next, we seek ν^* such that $\sum_{k=1}^n \theta_k^* = 1$, or alternatively, that $e^{c^* \sum_{k=1}^n \theta_k^*} = e^{c^*}$. Using Eq. S8, we can write

$$\prod_{k=1}^n \max \left\{ 1, (1 - \gamma_k^*) \frac{\nu^* - c^* \sum_{i=1}^{k-1} X_i}{\nu^* - c^* \sum_{i=1}^k X_i} \right\} = e^{c^* \sum_{k=1}^n \theta_k^*} = e^{c^*}. \quad [\text{S9}]$$

The product in Eq. S9 forms a piecewise continuous and monotonously decreasing function of ν^* , with breakpoints at the finite α_k 's. To observe these properties, we assume *w.l.o.g.* that all α_k 's

are finite and inspect the function's behavior when starting at $\nu_{ub}^* = \max_{1 \leq k \leq n} \{\alpha_k\}$ and then gradually decreasing ν^* . Clearly, the product yields 1 at ν_{ub}^* , but once $\nu^* < \nu_{ub}^*$, it increases due to at least one site k_0 for which $\nu^* < \alpha_{k_0}$ and which contributes

$$(1 - \gamma_k^*) \frac{\nu^{*-c^*} \sum_{i=1}^{k_0-1} X_i}{\nu^{*-c^*} \sum_{i=1}^{k_0} X_i} > 1$$

to the total product. As we continue decreasing ν^* , additional terms increase their contribution to more than 1, with transitions taking place at the α_k 's. Moreover, each of these "flexible" terms forms a strictly monotonously decreasing function of ν^* (once it exceeds 1) and so it increases its contribution as ν^* decreases. Importantly, each increasing term reflects a continuous increase in its corresponding θ . This way, the growth of the total product facilitates the expansion of certain θ_k^* 's while keeping others fixed at 0. Once $\sum_{k=1}^n \theta_k^* = 1$, Eq. S9 holds and a unique assignment for Θ^* is determined as follows: Suppose the intersection occurs at $\alpha_l \leq \nu_{opt}^* < \alpha_m$, then

$\theta_k^* = 0$ for all k such that $\alpha_k \leq \alpha_l$, and $\theta_k^* = \frac{1}{c^*} [\log \frac{\nu_{opt}^{*-c^*} \sum_{i=1}^{k-1} X_i}{\nu_{opt}^{*-c^*} \sum_{i=1}^k X_i} + \log(1 - \gamma_k^*)]$ otherwise. By replacing $\sum_{i=1}^k X_i$ with $\sum_{i=1}^{n+1} X_i - \sum_{i=k+1}^{n+1} X_i$ we obtain the expression for θ_k^* in Theorem 3 in the main text. At this point, we wish to note that this solution technique is known as *water filling* (1, 2), and thus we call the product in Eq. S9 the *water-filling function*. The technique is illustrated in Fig. S1B and is visualized as the flooding of a region with varying surface levels up to a constant amount of water.

Finally, we argued earlier that the domain $\nu^* \leq c^* \sum_{i=1}^k X_i$ is irrelevant to the solution. We also mentioned that the formula for the reactivities that is computed by Algorithm 1 forms a special case of Eq. S6 when $\nu^* = c^* \sum_{i=1}^{n+1} X_i$. Now, if all estimated reactivities are nonnegative when $\nu^* = c^* \sum_{i=1}^{n+1} X_i$, then ν^* is the desired threshold. However, if some reactivities are negative, then the sum of the positive ones must exceed 1, and hence ν^* should be increased to more than $c^* \sum_{i=1}^{n+1} X_i$ until the sum reaches 1. This implies that $\nu^* \geq c^* \sum_{i=1}^{n+1} X_i > c^* \sum_{i=1}^k X_i \forall k$ and also justifies setting $\theta_k^* = 0$ whenever $X_k = 0$, because Eq. S3 now takes the form $c^* \sum_{i=1}^{k-1} X_i + \lambda_k^* = \nu^*$ and must hold with $\lambda_k^* > 0$.

The *second* subproblem, hereafter called **P2**, entails the maximization of $\log \mathcal{L}(\Theta^*, \Gamma, c^*)$ when both Θ^* and c^* are kept fixed and when Γ lies in the unit hypercube $[0, 1]^n$. **P2**'s strong convexity is straightforward, as shown below.

Lemma 2. *P2 is a convex optimization problem.*

Proof: From the expression for the likelihood function in the main text we have

$$\begin{aligned} \log \mathcal{L}(\Theta^*, \Gamma, c^*) &= \sum_{k=1}^n Y_k \left[\sum_{i=1}^{k-1} \log(1 - \gamma_i) + \log \gamma_k \right] \\ &+ \sum_{k=1}^n X_k \left[c^* \left(\sum_{i=k+1}^n \theta_i^* - 1 \right) + \sum_{i=1}^{k-1} \log(1 - \gamma_i) \right. \\ &+ \left. \log(e^{c^* \theta_k^*} - 1 + \gamma_k) \right] + (X_{n+1} + Y_{n+1}) \\ &\times \sum_{i=1}^n \log(1 - \gamma_i) - c^* X_{n+1}, \end{aligned} \quad \text{[S10]}$$

which consists of constants and of a positively weighted sum of logarithms of affine functions of Γ . It is well known that convexity is preserved under affine transformation and that logarithms are strictly concave (1), thus establishing the log-likelihood's concavity. This fact, together with the imposed box constraints, completes the proof.

We are now in a position to prove Theorem 2 in the main text.

Proof of Theorem 2: A key observation here is that **P2** can be transformed into n independent *unconstrained* optimization problems, each in one variable. This is because $\log \mathcal{L}(\Theta^*, \Gamma, c^*)$ is separable in $\gamma_1, \dots, \gamma_n$ as follows:

$$\begin{aligned} \log \mathcal{L}(\Theta^*, \Gamma, c^*) &= \sum_{k=1}^n \left[\left(\sum_{i=k+1}^{n+1} (X_i + Y_i) \right) \log(1 - \gamma_k) + Y_k \log \gamma_k \right. \\ &+ \left. X_k \log(e^{c^* \theta_k^*} - 1 + \gamma_k) \right] + C, \end{aligned} \quad \text{[S11]}$$

where C is a constant. For simplicity, we introduce the constants $S_k = \sum_{i=k+1}^{n+1} (X_i + Y_i)$ and functions

$$l_k(\gamma_k) = S_k \log(1 - \gamma_k) + Y_k \log \gamma_k + X_k \log(e^{c^* \theta_k^*} - 1 + \gamma_k) \quad \text{[S12]}$$

for $1 \leq k \leq n$, such that $\log \mathcal{L}(\Theta^*, \Gamma, c^*) = C + \sum_{k=1}^n l_k(\gamma_k)$. Now, not only each function can be optimized separately, but the box constraints can be removed as well, as long as $Y_k > 0$ or $X_k > 0$ jointly with $\theta_k^* = 0$. To see this, note that $S_k > Y_{n+1} > 0$ and hence $\log(1 - \gamma_k)$ serves as a natural barrier that upper bounds $l_k(\gamma_k)$'s domain at 1. Similarly, $Y_k > 0$ (or $X_k > 0$, $\theta_k^* = 0$) imposes zero as a lower bound via $\log \gamma_k$. Optimization in the absence of the $\log \gamma_k$ term is also straightforward, as explained below.

Assume *w.l.o.g.* that $X_k, Y_k, \theta_k^* > 0$ and consider setting

$$l'_k(\gamma_k) = -\frac{S_k}{1 - \gamma_k} + \frac{Y_k}{\gamma_k} + \frac{X_k}{e^{c^* \theta_k^*} - 1 + \gamma_k} = 0, \quad \text{[S13]}$$

which reduces to a quadratic equation when $\gamma_k \notin \{0, 1, 1 - e^{c^* \theta_k^*}\}$ (recall that $e^{c^* \theta_k^*} > 1$). $l'_k(\gamma_k)$'s concavity asserts that the equation's solution (if it exists) corresponds to its global maximum. Now, $l'_k(\gamma_k)$ is well-defined and continuous over $(0, 1)$ and approaches ∞ near 0 and $-\infty$ near 1. Therefore, it must cross zero inside the unit interval at

$$\begin{aligned} \gamma_k^* &= \frac{1}{2(X_k + Y_k + S_k)} \left[X_k + Y_k - T_k(S_k + Y_k) \right. \\ &+ \left. \sqrt{[X_k + Y_k - T_k(S_k + Y_k)]^2 + 4T_k Y_k (X_k + Y_k + S_k)} \right], \end{aligned} \quad \text{[S14]}$$

where $T_k = e^{c^* \theta_k^*} - 1 > 0$ and where

$$X_k + Y_k + S_k = \sum_{i=k}^{n+1} (X_i + Y_i) = S_{k-1} \quad \text{[S15]}$$

stands for the total count of fragments that contain nucleotide k . When $Y_k = 0$ but $\theta_k^* > 0$, the solution simplifies to $\gamma_k^* = \frac{X_k - T_k S_k}{X_k + S_k}$, but $l'_k(\gamma_k)$'s domain is now extended into \mathbb{R}^- with $l'_k(\gamma_k) \rightarrow \infty$ near $-T_k$. For this reason, we set $\gamma_k^* = \max\{0, \frac{X_k - T_k S_k}{X_k + S_k}\}$ in this case. The last case of interest is when $\theta_k^* = 0$ (i.e., $T_k = 0$), and then $\gamma_k^* = \frac{X_k + Y_k}{X_k + Y_k + S_k} > 0$, which generalizes the initial estimate $\hat{\gamma}_k = \frac{Y_k}{\sum_{i=k}^{n+1} Y_i}$ to data that were aggregated from both channels.

The *first* optimization problem entails setting $(\Theta, \Gamma) = (\Theta^*, \Gamma^*)$ and seeking a positive c that maximizes $\log \mathcal{L}(\Theta^*, \Gamma^*, c)$. We refer to this problem as **P1** and establish its strong convexity next.

Lemma 3. *P1 is a convex optimization problem.*

Proof: The concavity of $\log \mathcal{L}(\Theta^*, \Gamma^*, c)$ with respect to c can be seen by following the same derivation as in Lemma 1. Specifically, it suffices to maximize

$$\mathcal{H}(c) = \mathcal{F}(\Theta^*, \Gamma^*, c) - cX_{n+1}, \quad \text{[S16]}$$

where $\mathcal{F}(\Theta^*, \Gamma^*, c)$ stands for $\mathcal{F}(\Theta = \Theta^*)$, with $c = c^*$ (see Eq. S1). One can readily see that $\mathcal{H}(c)$ consists of a sum of a linear function of c and of the functions $h_k(c) = \log(e^{\theta_k^*} - (1 - \gamma_k^*))$ for all k for which $\theta_k^* \neq 0$. Because each of the latter functions is symmetric with respect to θ_k and to c , their second-order derivatives with respect to c are all negative (see Lemma 1), and so is the derivative of their positively weighted sum. This establishes $\mathcal{H}(c)$'s concavity. At this point, it is worth noting that $\mathcal{H}(c)$'s domain takes the form $\mathcal{D}_{\mathcal{H}} = \{c | c > c_{lb}\}$, where $c_{lb} = \max_k: \theta_k^* \neq 0 \{ \frac{1}{\theta_k^*} \log(1 - \gamma_k^*) \}$. We remind the reader that $\log \mathcal{L}(\Theta, \Gamma, c)$ with $Y_{n+1} > 0$ constrains the γ_k 's to $0 < 1 - \gamma_k \leq 1$, and therefore $c_{lb} < 0$ if and only if $\gamma_k^* > 0$ for each k for which $\theta_k^* \neq 0$. In addition, because the logarithmic terms are not well-defined whenever $\theta_k^* = \gamma_k^* = 0$, the function pertains only to those positions where $\theta_k^* \neq 0$ or $\gamma_k^* \neq 0$. Importantly, we do not expect to encounter positions where $\theta_k^* = \gamma_k^* = 0$. This is because when $\theta_k^* = 0$, the k th term in problem P2 becomes $\log \gamma_k^*$, thus ensuring $\gamma_k^* > 0$ whenever $X_k \neq 0$. If $X_k = 0$, then we have $Y_k \neq 0$ (recall that all double-zero positions are excluded from optimization), in which case the term $Y_k \log \gamma_k^*$ serves the same purpose. Finally, note that the constraint $c > 0$ does not comply with the standard (stricter) definition of a convex optimization problem, which involves weak inequalities, i.e., $c \geq 0$ (1). Yet, this subtlety is accounted for in the next lemma, and we refer to P1 as a convex optimization problem in the broader sense.

Problem P1 does not pose a significant computational challenge as it involves one variable. Yet, its solution can be further simplified by relaxing the constraint $c > 0$. This gives rise to an unconstrained optimization problem, where a concave $\mathcal{H}(c)$ is maximized over the semibounded interval $\mathcal{D}_{\mathcal{H}}$. Because $\mathcal{D}_{\mathcal{H}}$ may contain negative values, the relaxation might result in a negative maximizing argument. The following lemma shows that a maximum for the unconstrained problem is indeed attained within its semibounded domain and that it can be readily used to infer P1's solution.

Lemma 4. $\mathcal{H}(c)$ attains a unique maximum over its domain $\mathcal{D}_{\mathcal{H}} = \{c | c > c_{lb}\}$. If $\arg \max \mathcal{H}(c) \leq 0$, then problem P1 is infeasible. Otherwise, $\arg \max \mathcal{H}(c)$ is also the solution of P1.

Proof: We first demonstrate that $\mathcal{H}(c)$ attains a maximum over $\mathcal{D}_{\mathcal{H}}$ by showing that its derivative crosses zero inside the domain. The derivative takes the form

$$\mathcal{H}'(c) = \sum_{k=1}^n X_k \left[\sum_{l=k+1}^n \theta_l^* + \frac{\theta_k^*}{1 - (1 - \gamma_k^*)e^{-c\theta_k^*}} - 1 \right] - X_{n+1} \quad \text{[S17]}$$

and is well-defined and continuous over $\mathcal{D}_{\mathcal{H}}$. When c approaches c_{lb} from the right, we have $\frac{\theta_k^*}{1 - (1 - \gamma_k^*)e^{-c\theta_k^*}} \rightarrow \infty$ for at least one position k with $\theta_k^* \neq 0$, implying that $\mathcal{H}'(c) > 0$ for sufficiently small c values. In contrast, when $c \rightarrow \infty$, we have $\frac{\theta_k^*}{1 - (1 - \gamma_k^*)e^{-c\theta_k^*}} \rightarrow \theta_k^*$ for each k with $\theta_k^* > 0$ (and a zero-term otherwise), and therefore $\mathcal{H}'(c) \rightarrow \sum_{k=1}^n X_k (\sum_{l=k}^n \theta_l^* - 1) - X_{n+1}$. The expression on the right-hand side must be negative, because $\sum_{l=k}^n \theta_l^* \leq 1$ with strict inequality for some k 's and because we expect to have $X_{n+1} > 0$ (see *Modeling Chemical Modification* in the main text for reasoning). It now follows from the intermediate value theorem that $\mathcal{H}'(c^*) = 0$ for some $c^* \in \mathcal{D}_{\mathcal{H}}$, which is where the unique maximum is located. Finally, $\text{sign}(c^*)$ determines whether c^* solves

P1 as well. In case it is negative, we have $\mathcal{H}(c) < \mathcal{H}(0)$ for all $c > 0$, and because $\mathcal{H}(c)$ monotonously decreases over \mathbb{R}^+ P1 is infeasible. Otherwise, c^* maximizes $\mathcal{H}(c)$ over \mathbb{R}^+ and is P1's sought solution.

Concavity of the log-Likelihood in the Entire Parameter Space.

Although the decomposition of the global optimization problem into local optimization steps lends itself to an efficient semianalytic solution, it is also of interest to investigate the function's concavity with respect to the entire coordinate system. Such global concavity has two implications: first, every local maximum is also global, and thus Algorithm 2 should converge to a global maximum even when initialized at a single starting point. Second, it facilitates the solution of the entire optimization problem using numerical interior-point methods. Yet, it is not obvious that the latter approach is advantageous over Algorithm 2, as it might be more computationally intensive and/or yield less accurate solutions. Hence, even if the function is globally concave, the performance of interior-point methods needs to be carefully examined and compared to that of Algorithm 2's.

Before we study $\log \mathcal{L}(\Theta, \Gamma, c)$ as a function of all parameters, we set $c = c^*$ and investigate its concavity with respect to (Θ, Γ) . We can further simplify analysis by following the lines of Theorem 2's proof and observing that $\log \mathcal{L}(\Theta, \Gamma, c^*)$ is separable in the pairwise variables (θ_k, γ_k) as follows:

$$\log \mathcal{L}(\Theta, \Gamma, c^*) = \sum_{k=1}^n l_k(\theta_k, \gamma_k) + L(\Theta) + C. \quad \text{[S18]}$$

Here, $L(\Theta)$ represents a linear function of Θ , whereas C and $l_k(\theta_k, \gamma_k)$ pertain to the constant and functions that were introduced in Theorem 2's proof, respectively. Specifically,

$$l_k(\theta_k, \gamma_k) = S_k \log(1 - \gamma_k) + Y_k \log \gamma_k + X_k \log(e^{c^* \theta_k} - 1 + \gamma_k), \quad \text{[S19]}$$

where $S_k = \sum_{i=k+1}^{n+1} (X_i + Y_i)$. Because $\log \mathcal{L}(\Theta, \Gamma, c^*)$ is concave in (Θ, Γ) if and only if $l_k(\theta_k, \gamma_k)$ is concave in (θ_k, γ_k) for all k , it suffices to investigate the properties of the two-dimensional function $l_k(\theta_k, \gamma_k)$. We start by showing that its right-hand term, that is,

$$f_k(\theta_k, \gamma_k) = \log(e^{c^* \theta_k} - (1 - \gamma_k)), \quad \text{[S20]}$$

is *not* concave in (θ_k, γ_k) . Before we proceed, we stress that $f_k(\theta_k, \gamma_k)$'s domain, consisting of all (θ_k, γ_k) such that $\gamma_k > 1 - e^{c^* \theta_k}$, does *not* form a convex set. This can be seen from the fact that the domain consists of the epigraph of the strictly concave function $\gamma_k(\theta_k) = 1 - e^{c^* \theta_k}$. In light of the importance of a domain's convexity in determining a function's concavity (1), we artificially restrict the function's domain to the following convex subset of the original domain, which contains the entire feasible parameter space:

$$\mathcal{D}_{\mathcal{S}} = \{(\theta_k, \gamma_k): 0 \leq \theta_k \leq 1, 0 \leq \gamma_k \leq 1\}.$$

Lemma 5. The function $f_k: \mathbb{R}^2 \rightarrow \mathbb{R}$, with domain $\mathcal{D}_{\mathcal{S}} = \{(\theta_k, \gamma_k): 0 \leq \theta_k \leq 1, 0 \leq \gamma_k \leq 1\}$ and given by $f_k(\theta_k, \gamma_k) = \log(e^{c^* \theta_k} - (1 - \gamma_k))$, is *not* concave in (θ_k, γ_k) .

Proof: Because $f_k(\theta_k, \gamma_k)$ is continuous and twice differentiable in its domain, its concavity can be inferred from the negative semidefiniteness of its Hessian on the interior of $\mathcal{D}_{\mathcal{S}}$. The Hessian is given by

$$H_{f_k} = -[e^{c^*\theta_k} - (1 - \gamma_k)]^{-2} \begin{bmatrix} c^{*2}e^{c^*\theta_k}(1 - \gamma_k) & c^*e^{c^*\theta_k} \\ c^*e^{c^*\theta_k} & 1 \end{bmatrix}, \quad [\text{S21}]$$

and when multiplied by $V = (x, y)$ from both sides we obtain

$$VH_{f_k}V^T = -[e^{c^*\theta_k} - (1 - \gamma_k)]^{-2}[c^{*2}e^{c^*\theta_k}(1 - \gamma_k)x^2 + y^2 + 2c^*e^{c^*\theta_k}xy] \quad [\text{S22}]$$

$$= -[e^{c^*\theta_k} - (1 - \gamma_k)]^{-2}[(c^*e^{c^*\theta_k}x + y)^2 + c^{*2}e^{c^*\theta_k}(1 - \gamma_k - e^{c^*\theta_k})x^2]. \quad [\text{S23}]$$

Consider (x, y) such that $x \neq 0$ and $c^*e^{c^*\theta_k}x + y = 0$, in which case

$$VH_{f_k}V^T = -[e^{c^*\theta_k} - (1 - \gamma_k)]^{-2}c^{*2}e^{c^*\theta_k}(1 - \gamma_k - e^{c^*\theta_k})x^2 \quad [\text{S24}]$$

$$= [e^{c^*\theta_k} - (1 - \gamma_k)]^{-1}c^{*2}e^{c^*\theta_k}x^2 > 0, \quad [\text{S25}]$$

where the positive sign follows from the fact that $f_k(\theta_k, \gamma_k)$ is defined only when $e^{c^*\theta_k} - (1 - \gamma_k) > 0$ and because $c^* > 0$. Eq. S25 implies that H_{f_k} is not negative semidefinite and hence $f_k(\theta_k, \gamma_k)$ is not concave. Note that H_{f_k} is not positive semidefinite as well, because a choice of $V = (0, y)$ with $y > 0$ yields a negative product in Eq. S22. Therefore, $f_k(\theta_k, \gamma_k)$ is neither convex nor concave on the domain of interest.

Importantly, the fact that $f_k(\theta_k, \gamma_k)$ is not concave does not imply that $l_k(\theta_k, \gamma_k) = S_k \log(1 - \gamma_k) + Y_k \log \gamma_k + X_k f_k(\theta_k, \gamma_k)$ is not concave as well, because the first two terms are strictly concave functions. In what follows, we attempt to investigate $l_k(\theta_k, \gamma_k)$'s concavity by examining its Hessian, which takes the form

$$H_{l_k} = - \begin{bmatrix} \frac{X_k c^{*2} e^{c^*\theta_k} (1 - \gamma_k)}{[e^{c^*\theta_k} - (1 - \gamma_k)]^2} & \frac{X_k c^* e^{c^*\theta_k}}{[e^{c^*\theta_k} - (1 - \gamma_k)]^2} \\ \frac{X_k c^* e^{c^*\theta_k}}{[e^{c^*\theta_k} - (1 - \gamma_k)]^2} & M \end{bmatrix}, \quad [\text{S26}]$$

where $M = \frac{X_k}{[e^{c^*\theta_k} - (1 - \gamma_k)]^2} + \frac{S_k}{(1 - \gamma_k)^2} + \frac{Y_k}{\gamma_k^2}$. We now exploit the fact that H_{l_k} is two-dimensional and hence its determinant's sign is indicative of its definiteness. If the sign is negative, then the two eigenvalues must be of different signs and the matrix is neither positive nor negative semidefinite, or equivalently, the function is neither concave nor convex. The determinant is given by

$$\det H_{l_k} = \frac{X_k c^{*2} e^{c^*\theta_k}}{[e^{c^*\theta_k} - (1 - \gamma_k)]^2} \left[\frac{-X_k}{e^{c^*\theta_k} - (1 - \gamma_k)} + \frac{S_k}{1 - \gamma_k} + \frac{Y_k(1 - \gamma_k)}{\gamma_k^2} \right], \quad [\text{S27}]$$

and if $X_k > 0$ then

$$\text{sign}(\det H_{l_k}) = \text{sign} \left(\frac{-X_k}{e^{c^*\theta_k} - (1 - \gamma_k)} + \frac{S_k}{1 - \gamma_k} + \frac{Y_k(1 - \gamma_k)}{\gamma_k^2} \right). \quad [\text{S28}]$$

Recall that $S_k > 0$ and let us assume for simplicity that X_k and Y_k are also positive. Recall also that we are concerned with all the points in $\mathcal{D}\mathcal{S}$'s interior where $0 < \theta_k < 1$ and $0 < \gamma_k < 1$. We thus wish to explore the existence of a feasible point for which the expression in Eq. S28 assumes a negative value. Although this expression reduces to a cubic equation in γ_k , finding its real root analytically is intractable. Instead, we observe that $\frac{X_k}{e^{c^*\theta_k} - (1 - \gamma_k)}$ depends on both θ_k and γ_k , whereas the other two terms depend

only on γ_k . Additionally, this term increases in magnitude as $\gamma_k \rightarrow 0$ and as $e^{c^*\theta_k} \rightarrow 1$. However, we also have $\frac{Y_k(1 - \gamma_k)}{\gamma_k^2} \rightarrow \infty$ when $\gamma_k \rightarrow 0$ and $\frac{S_k}{1 - \gamma_k} \rightarrow \infty$ when $\gamma_k \rightarrow 1$. This means that for any given $c^*\theta_k$ value, $\frac{X_k}{e^{c^*\theta_k} - (1 - \gamma_k)}$ needs to offset an increasing (and unbounded) positive term at both limits of the interval $0 < \gamma_k < 1$. It is therefore not obvious whether there exist θ_k and γ_k for which $\frac{-X_k}{e^{c^*\theta_k} - (1 - \gamma_k)} + \frac{S_k}{1 - \gamma_k} + \frac{Y_k(1 - \gamma_k)}{\gamma_k^2} < 0$.

Instead, we resort to considering the limiting case where $\theta_k = 0$ (on $\mathcal{D}\mathcal{S}$'s boundary), because this is where $\frac{X_k}{e^{c^*\theta_k} - (1 - \gamma_k)}$ is maximized for each γ_k . Here, we can test whether the simplified condition

$$\frac{-X_k}{\gamma_k} + \frac{S_k}{1 - \gamma_k} + \frac{Y_k(1 - \gamma_k)}{\gamma_k^2} < 0 \quad [\text{S29}]$$

holds for some $0 < \gamma_k < 1$. By converting the latter condition into the equivalent inequality

$$(S_k + Y_k + X_k)\gamma_k^2 - (2Y_k + X_k)\gamma_k + Y_k < 0, \quad [\text{S30}]$$

it is easy to show that it is satisfied within some subinterval of $(0, 1)$, provided that

$$X_k^2 > 4Y_k S_k. \quad [\text{S31}]$$

This is because its left-hand side is positive at $\gamma_k \in \{0, 1\}$ and its minimum is attained at $0 < \tilde{\gamma}_k = \frac{2Y_k + X_k}{2(S_k + Y_k + X_k)} < 1$. Now, if the expression's discriminant is positive, it must cross zero twice inside $(0, 1)$ and attain negative values between the two intersection points. Finally, constraining the discriminant in this manner yields inequality S31. When S31 is satisfied, we can choose a sufficiently small $c^*\theta_k$ such that $e^{c^*\theta_k} - 1 < \epsilon$ and $\frac{X_k}{e^{c^*\theta_k} - (1 - \gamma_k)} \approx \frac{X_k}{\gamma_k}$. This, in turn, guarantees that $\frac{-X_k}{e^{c^*\theta_k} - (1 - \gamma_k)} + \frac{S_k}{1 - \gamma_k} + \frac{Y_k(1 - \gamma_k)}{\gamma_k^2} < 0$ for the chosen values.

When S31 does not hold, it follows that $\frac{X_k}{e^{c^*\theta_k} - (1 - \gamma_k)} < \frac{X_k}{\gamma_k} \leq \frac{S_k}{1 - \gamma_k} + \frac{Y_k(1 - \gamma_k)}{\gamma_k^2}$, or that $\text{sign}(\det H_{l_k}) > 0$. Notably, a positive determinant implies that $l_k(\theta_k, \gamma_k)$ is either convex or concave, but because $f_k(\theta_k, \gamma_k)$ is neither concave nor convex (see the proof of Lemma 5) and $\log(1 - \gamma_k)$ and $\log \gamma_k$ are strictly concave, $l_k(\theta_k, \gamma_k)$ must also be concave in this case. Finally, when S31 does not hold for every $1 \leq k \leq n$, we can conclude that $\log \mathcal{L}(\Theta, \Gamma, c^*)$ is concave in (Θ, Γ) .

Our analysis thus suggests that the global concavity of $\log \mathcal{L}(\Theta, \Gamma, c^*)$ is data-dependent, as captured by the condition in S31. In light of this condition, we consider the case where $Y_k = 0$ but $X_k > 0$, which, in our experience, is commonly observed, albeit for very few positions. For such positions, S31 clearly holds, thereby eliminating the log-likelihood's global concavity. Although one may arbitrarily set $Y_k = 1$ for such positions, S31 may or may not be satisfied following such correction, depending on the corresponding values of X_k and S_k .

As final remarks, our analysis considered the domain $\mathcal{D}\mathcal{S} = \{(\theta_k, \gamma_k) : 0 \leq \theta_k \leq 1, 0 \leq \gamma_k \leq 1\}$, but, in fact, optimization takes place over a convex subset of this domain, in which $\sum_{k=1}^n \theta_k = 1$. It is possible that under a fixed-sum restriction, $\log \mathcal{L}(\Theta, \Gamma, c^*)$ exhibits global concavity. However, analysis of the type we conducted above does not apply in this case due to two reasons: (i) only directions $V = (v_1, \dots, v_n)$ such that $\sum_{k=1}^n v_k = 0$ are feasible under the additional restriction, and (ii) one cannot leverage on the separability of $\log \mathcal{L}(\Theta, \Gamma, c^*)$ into two-dimensional functions, as these are now linked by the fixed-sum constraint. Therefore, such analysis exceeds the scope of this work and is a topic for

sidered, but they consistently yielded 100% hits, which is indicative of overall low sensitivity to the approximation and is also the reason why we decreased the resolution of our accuracy measure.

We tested a variety of model scenarios. First, we assumed that Θ is uniform and changed the Poisson rate within the range $c = 0.25$ – 2 . We also considered two uniform distributions: one with $n = 10$ reactive sites and one with $n = 50$ reactive sites. Table S1 shows that accuracy improves with decreasing modification rate and with increasing molecule size (more precisely, with increasing number of reactive sites). This is expected, because, in general, the statistics of sampling with and without replacement become similar when the number of sites to draw from is significantly larger than the number of draws (i.e., modifications). Second, for a uniform Θ , we varied n while keeping the rate fixed at $c = 1$ and at $c = 2$. One can see the same trend, where increased lengths lead to lower sensitivity, and that high quality is achieved for molecules with 40 or more reactive sites, a scenario which, in our experience, is realistic. Third, we considered four other Θ distributions, as follows: A decreasing exponential and an increasing exponential were chosen as extreme examples of unbalanced distributions, although we do not think they are representative of realistic reactivity profiles as these typically display some form of symmetry due to base pairing. A decreasing exponential of length $n = 50$ takes the form $a \times 0.2 \times (0.9)^i$ (a is a normalization constant), starting at approximately 0.1 and reaching very close to zero at $i = 50$, where the increasing exponential is its flipped version. We also generated two profiles that reflect our understanding of the crude form of a single hairpin ($n = 20$) and a concatenation of two hairpins ($n = 50$). It can be seen from Table S1 that the model had low sensitivity to the approximation under the first three distributions, whereas estimation was less accurate in the single hairpin case.

Robustness of estimates under different model distributions. We tested estimation quality under different model parameters by simulating a total of 5 million draws from the induced fragment length distributions in the (+) and (–) channels, and then processing the counts with Algorithms 1 and 2. We repeated this process 500 times. Here, we considered the average fraction of hits at resolutions of 15%, 10%, and 5% of the true Θ as well as the standard deviation around the average (reported in Table S2). We also computed the standard deviation of the estimated reactivity per site. The standard deviation (per site and per interval) was negligible in all cases considered. We first assumed that $\gamma_k = 0.01$ for all k and simulated the model for the five Θ distributions mentioned in the previous subsection as well as for a uniform, a decreasing exponential, and an increasing exponential, all of size $n = 10$. Subsequently, we relaxed the fixed- Γ assumption and

added spikes of magnitude 0.15 on top of the $\gamma_k = 0.01$ background. The spike magnitude was chosen based on spikes we observed in the control experiments of the RNase P and the pT181 molecules. We tested the effects of two and four such spikes (distributed evenly across the molecule) when Θ is uniform and when Θ represents two concatenated hairpins. The results are summarized in Table S2.

Read Alignment Details. Reads for reverse transcriptase (RT) fragments were first split into 1M7-treated and untreated pools by examining the 4 nucleotide handle sequence on the 5' end of the read generated from the 3' end of each RNA probed in the experiment. Reads with an RRRY handle identified (+) fragments and those with YYR identified (–) fragments. This handle was then trimmed from each read to allow alignment of the reads to probed RNAs. Reads were then trimmed for A_adapter_b and A_adapter_t using the FASTX toolkit, because RT products shorter than the length of a sequencing read will produce reads with adapter at their 3' ends.

Paired reads were optimally aligned to the probed RNAs using Bowtie 0.12.8 (3) with the parameters (–best X 2000 y), and allowing no mismatches in the alignments with parameters (–v 0). Bowtie suppressed alignments for reads that mapped ambiguously to the probed RNAs (–m 1), though because the bar code in each read identified each read with its target RNA and alignments were required to be perfect, this step was not strictly necessary. The 3' end of each fragment alignment (toward the 5' end of the probed RNA) corresponds to the point at which RT stopped. Two counters (one for the 1M7 condition and one for the control condition) at each position of each probed RNA position were used to track RT stopping points. The (+) counter for a probed RNA at position i was incremented when a fragment from the (+) pool aligned to the RNA starting at position $(i + 1)$ and ending at the RNAs 3' end. The (–) counter was incremented for (–) pool fragment alignments. These RT-stop counts were then used to calculate maximum likelihood reactivities.

Source Code. Source code to implement read mapping and ML estimation can be found at <http://bio.math.berkeley.edu/SHAPE-Seq/> as either a .tar.gz or .zip file. After unpacking with the command `tar -xzf Aviran_Trapnell_SHAPE-Seq_analysis_code_PNAS_2011.tar.gz` or `unzip Aviran_Trapnell_SHAPE-Seq_analysis_code_PNAS_2011.zip` a directory called `spats-0.0.1` will be created. Please follow instructions in the README file to compile and execute the read mapping code (spats). Matlab source code to implement the ML estimation can be found in `spats-0.0.1/src/matlab`.

1. Boyd S, Vandenberghe L (2004) *Convex Optimization*. (Cambridge Univ Press, Cambridge, UK).
2. Cover TM, Thomas JA (2006) *Elements of Information Theory*. (John Wiley, Hoboken, NJ), 2nd Ed.

3. Langmead B, Trapnell C, Pop M, Salzberg SL (2009) Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 10:R25.

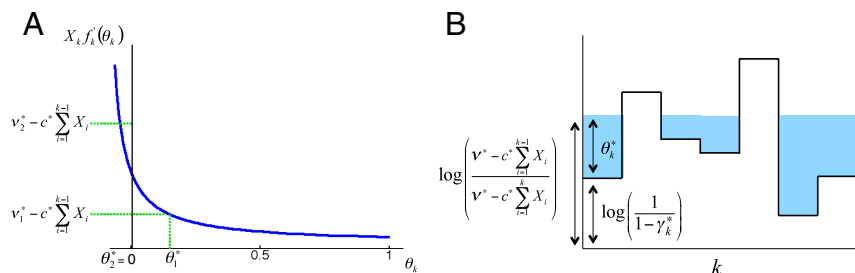


Fig. S1. Illustration of the use of water filling to optimize Θ . (A) Illustration of the function $f_k(\theta_k)$ and of the solution obtained for a single KKT condition under two instances of ν^* . In the first instance (ν_1^*), a strictly positive solution is obtained, whereas in the second instance (ν_2^*), a zero reactivity is assigned. (B) Illustration of the water-filling technique following ref. 1. The height of each bar is computed from the corresponding parameter γ_k^* . The gradual decreasing of ν^* can be viewed as flooding the bars up to a level of $\log(\nu^* - c^* \sum_{i=1}^{k-1} X_i) - \log(\nu^* - c^* \sum_{i=1}^k X_i)$, until the total quantity of water equals c^* . The level of water above each bar (in light blue) corresponds to the optimal θ_k^* .

1. Boyd S, Vandenberghe L (2004) *Convex Optimization* (Cambridge Univ Press, Cambridge, UK).

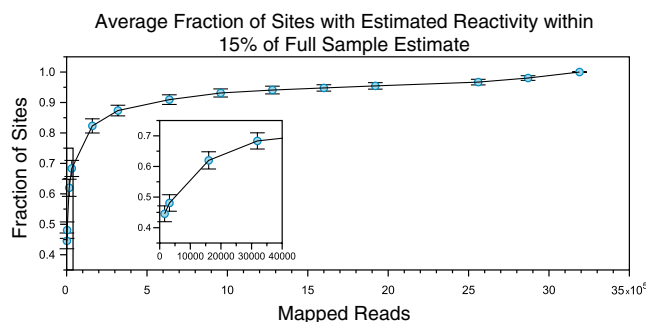


Fig. S2. Dependence of the ML estimates on the number of mapped reads for the pT181 molecule. Varying numbers of reads were drawn at random from the observed dataset and used to generate the ML reactivity estimate for each nucleotide. The fraction of these reactivities that were within 15% of the ML estimate of the full dataset was recorded. This was repeated 200 times for each specified number of reads. Points represent the average over these 200 draws and error bars represent standard deviations. Very accurate estimates are obtained with only 300,000 reads (10% of the full number of reads), and over half of the nucleotide positions show accurate reactivities with as little as 16,000 reads (0.5% of the full number of reads).

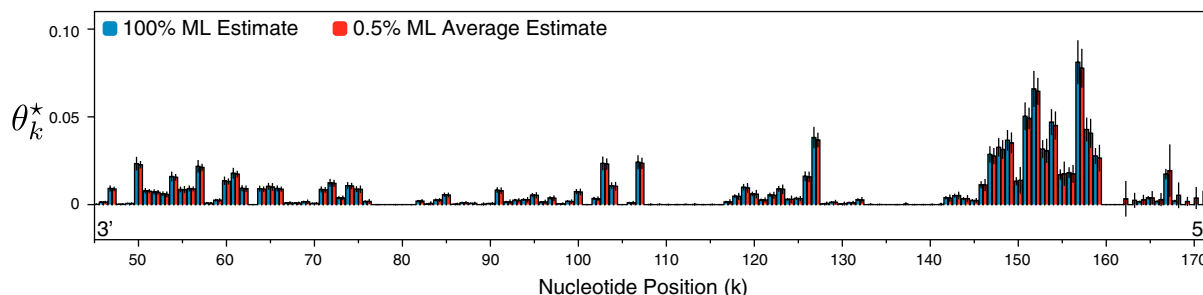


Fig. S3. Sample ML estimates based on a fraction of the observed data for the pT181 molecule. The full ML reactivity estimates for each nucleotide of the pT181 molecule (see Fig. 2 in the main text) were calculated from the 3.2 million mapped reads and are shown in blue with $\pm 15\%$ of the estimate plotted as error bars. The red bars represent an average of the ML reactivity estimate at each nucleotide, calculated from 200 independent random draws of 16,000 observed reads (0.5% of the total) with error bars representing the standard deviation of the estimate at each nucleotide.

Table S1. Model sensitivity with respect to the approximation of repeated modification per site

n	Θ	c	Average fraction of reactivities within $\pm 10\%$ of known value	Average fraction of reactivities within $\pm 5\%$ of known value	Average fraction of reactivities within $\pm 1\%$ of known value	
10	uniform	0.25	1	1	0.7	
		0.50	1	1	0.4	
		1.00	1	0.9	0.2	
		1.50	1	0.6	0.1	
		2.00	0.9	0.4	0.1	
50		0.25	1	1	1	
		0.50	1	1	1	
		1.00	1	1	0.86	
		1.50	1	1	0.72	
		2.00	1	1	0.44	
20		1	1	1	0.35	
40			1	1	0.83	
60			1	1	0.98	
80			1	1	0.99	
20	one hairpin	2	1	0.9	0.2	
40			1	1	0.35	
60			1	1	0.48	
80			1	1	0.61	
20			1	0.85	0.35	
50	uniform	1	1	0.55	0.15	
			decreasing	1	1	0.92
			exponential			
			increasing	1	1	0.06
			exponential	1	1	0.54
two hairpins						

Table S2. Robustness of estimates under different model distributions

n	Θ	Γ	Fraction of reactivities within $\pm 15\%$ of known value		Fraction of reactivities within $\pm 10\%$ of known value		Fraction of reactivities within $\pm 5\%$ of known value	
			Average	Standard deviation	Average	Standard deviation	Average	Standard deviation
			50	uniform	fixed (0.01)	1	0	1
	decreasing		0.91	0.03	0.86	0.04	0.76	0.04
	exponential							
	increasing		0.97	0.02	0.93	0.03	0.83	0.04
	two hairpins		0.99	0.01	0.97	0.02	0.86	0.04
10	uniform	fixed (0.01)	1	0	1	0	1	0
	decreasing		1	0	1	0	1	0
	exponential							
	increasing		1	0	1	0	1	0
	one hairpin		1	0	1	0	0.99	0.02
50	uniform	two spikes (0.01,0.15)	1	0	1	0	0.99	0.01
		four spikes (0.01,0.15)	1	0	1	0	0.99	0.01
	two hairpins	two spikes (0.01,0.15)	0.98	0.02	0.95	0.03	0.83	0.04
		four spikes (0.01,0.15)	0.96	0.02	0.92	0.03	0.81	0.04