

Precision of Readout at the hunchback Gene: Analyzing Short Transcription Time Traces in Living Fly Embryos

Jonathan Desponds, Huy Tran, Teresa Ferraro, Tanguy Lucas, Carmina Perez Romero, Aurelien Guillou, Cecile Fradin, Mathieu Coppey, Nathalie Dostatni, Aleksandra M. Walczak

► To cite this version:

Jonathan Desponds, Huy Tran, Teresa Ferraro, Tanguy Lucas, Carmina Perez Romero, et al.. Precision of Readout at the hunchback Gene: Analyzing Short Transcription Time Traces in Living Fly Embryos. PLoS Computational Biology, 2016, 12 (12), pp.e1005256. 10.1371/journal.pcbi.1005256 . hal-01470230

HAL Id: hal-01470230 https://hal.sorbonne-universite.fr/hal-01470230

Submitted on 17 Feb 2017 $\,$

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



Citation: Desponds J, Tran H, Ferraro T, Lucas T, Perez Romero C, Guillou A, et al. (2016) Precision of Readout at the *hunchback* Gene: Analyzing Short Transcription Time Traces in Living Fly Embryos. PLoS Comput Biol 12(12): e1005256. doi:10.1371/ journal.pcbi.1005256

Editor: Saurabh Sinha, University of Illinois at Urbana-Champaign, UNITED STATES

Received: July 13, 2016

Accepted: November 19, 2016

Published: December 12, 2016

Copyright: © 2016 Desponds et al. This is an open access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Additional data are available from http://xfer.curie.fr/get/ GmJzLUbF1JU/mov.zip.

Funding: This work was supported by a Marie Curie MCCIG grant No. 303561 (AMW), PSL IDEX REFLEX (ND, AMW, MC), ARC PJA20151203341 (ND), ANR-11-LABX-0044 DEEP Labex (ND), ANR-11-BSV2-0024 Axomorph (ND and AMW) and PSL ANR-10-IDEX-0001-02. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. RESEARCH ARTICLE

Precision of Readout at the *hunchback* Gene: Analyzing Short Transcription Time Traces in Living Fly Embryos

Jonathan Desponds^{1,2,3}, Huy Tran^{1,2,3}, Teresa Ferraro^{1,2,3}, Tanguy Lucas^{2,3,4}, Carmina Perez Romero⁵, Aurelien Guillou^{2,3,4}, Cecile Fradin⁵, Mathieu Coppey^{2,3,4}, Nathalie Dostatni^{2,3,4}, Aleksandra M. Walczak^{1,2,3,*}

1 Ecole Normale Superieure, PSL Research University, Paris, France, 2 Sorbonne Universités, UPMC Univ Paris 06, Paris, France, 3 UMR3664/UMR168/UMR8549, CNRS, Paris, France, 4 Institut Curie, PSL Research University, Paris, France, 5 McMaster University, Hamilton, Ontario, Canada

* awalczak@lpt.ens.fr

Abstract

The simultaneous expression of the *hunchback* gene in the numerous nuclei of the developing fly embryo gives us a unique opportunity to study how transcription is regulated in living organisms. A recently developed MS2-MCP technique for imaging nascent messenger RNA in living *Drosophila* embryos allows us to quantify the dynamics of the developmental transcription process. The initial measurement of the morphogens by the *hunchback* promoter takes place during very short cell cycles, not only giving each nucleus little time for a precise readout, but also resulting in short time traces of transcription. Additionally, the relationship between the measured signal and the promoter state depends on the molecular design of the reporting probe. We develop an analysis approach based on tailor made autocorrelation functions that overcomes the short trace problems and quantifies the dynamics of transcription initiation. Based on live imaging data, we identify signatures of bursty transcription initiation from the *hunchback* promoter. We show that the precision of the expression of the *hunchback* gene to measure its position along the anterior-posterior axis is low both at the boundary and in the anterior even at cycle 13, suggesting additional post-transcriptional averaging mechanisms to provide the precision observed in fixed embryos.

Author Summary

The fly embryo provides a natural laboratory to study the dynamics of transcription and its implications for the developing organism. Using live imaging experiments we investigate the nature of transcription regulation of the *hunchback* gene—the first to read out the maternal Bicoid gradient. While traditional time trace analysis methods based on OFF time distributions or autocorrelation functions fail for short signals, our tailored autocorrelation function overcomes these limitations revealing bursty dynamics that is reproducible between cell cycles and embryos. The inferred rates result in a lot of variability in the

Competing Interests: The authors have declared that no competing interests exist.

readout of nuclei sensing similar Bicoid concentrations, suggesting additional readout mechanisms than a one-to-one mapping of the input onto the output.

Introduction

During development the different identities of cells are determined by sequentially expressing particular subsets of genes in different parts of the embryo. Proper development relies on the correct spatial-temporal assignment of cell types. In the fly embryo, the initial information about the position along the anterior-posterior (AP) axis is encoded in the exponentially decaying Bicoid gradient. The simultaneous expression of the Bicoid target gene *hunchback* in the multiple nuclei of the developing fly embryo gives us a unique opportunity to study how transcription is regulated and controlled in a living organism [1, 2]. Despite many downstream points where possible mistakes can be corrected [1, 3, 4], the initial mRNA readout of the maternal Bicoid gradient by the *hunchback* gene is remarkably accurate and reproducible between embryos [5, 6]: it is highly expressed in the anterior part. This precision is even more surprising given the very short duration of the cell cycles (6–15 minutes) during which the initial Bicoid readout takes place and the intrinsic molecular noise in transcription regulation [7–9].

Even though most of our understanding of transcription regulation in the fly embryo comes from studies of fixed samples, gene expression is a dynamic process. The process involves the assembly of the transcription machinery and depends on the concentrations of the maternal gradients [10]. Recent studies based on single-cell temporal measurements of a short lived luciferase reporter gene under the control of a number of promoters in mouse fibroblast cell cultures [11, 12] and experiments in *E. Coli* and yeast populations [13–16] have quantitatively confirmed that mRNAs are generally produced in bursts, which result from periods of activation and inactivation. In early fly development, what are the dynamical properties of transcription initiation that allow for the concentration of the Bicoid gradient and other maternal factors to be measured in these short intervals between mitoses?

In order to quantitatively describe the events involved in transcription initiation, we need to have a signature of this process in the form of time dependent traces of RNA production. Recently, live imaging techniques have been developed to simultaneously track the RNA production in all nuclei throughout the developmental period from nuclear cycle 11 to cycle 14 [17, 18]. In these experiments, an MS2-binding cassette is placed directly under the control of an additional copy of a proximal *hunchback* promoter. As this reporter gene is transcribed, mRNA loops are expressed that bind fluorescent MCP proteins. Their accumulation at the transcribed locus gives an intense localized signal above the background level of unbound MCP proteins (Fig 1C) [19]. By monitoring the developing embryo, we obtain for each nucleus a time dependent fluorescence trace that is indicative of the dynamics of transcription regulation at the *hunchback* promoter (Fig 1B, 1D and 1F).

However the fluorescent time traces inevitably provide an indirect observation of the transcription dynamics. The signal is noisy, convoluting both experimental and intrinsic noise with the properties of the probe: the jitter in the signal is not necessary indicative of actual gene switching but could simply result from a momentarily decrease in the recording of the intensity. To obtain a sufficiently strong intensity of the signal to overcome background fluorescence, a long probe with a large number of loops is needed, which introduces a buffering time. In the current experiments the minimal buffering time is the time needed to transcribe a



Fig 1. Transcription dynamics in the fly embryo. (A) The three models of transcription dynamics considered in this paper. From left to right: the two state model, the cycle model and the Gamma model (see SI Sections B, D and E). (B) Example of the promoter state dynamics (either ON or OFF) as a function of time. We assume that the polymerase is abundant and every time the promoter is ON and is not flanked by the previous polymerase a new polymerase will start transcribing. The function X(t) in black is non-zero when a polymerase is occupying the transcription initiation site and zero otherwise. (C) In the ON state, the promoter (*Pr*) is accessible to RNA polymerases (Pol II) that initiate the transcription of the target gene and the 24× MS2 loops. As the 24× target mRNA is elongated MCP-GFP fluorescent molecules bind a detectable fluorescence signal. (D) MCP-GFP molecules labeling several mRNAs co-localize at the transcription loci, which appear as green spots under the confocal microscope. The spot intensities are then extracted over time and classified by each nuclei's position in the *Drosophila* embryo as Anterior, Boundary and Posterior. The spatial resolution of the spots is limited by the Abel limit, which is ~ 200*nm*. The ability to identify spots is also limited by the background level of free MCP-GFP. Typical spot sizes are ~ 260*nm*, giving an upper bound on the size of the transcription site. (E) The gene is divided into *r* sites of size 150 base pairs, indexed by *i*. The presence or absence of a polymerase. (F) A cartoon representing the type of experimental signal we analyze (see S1 Fig for real traces): one spot's intensity as a function of time, corresponding to the arrival of RNA polymerases in (E) and the promoter state in (B).

doi:10.1371/journal.pcbi.1005256.g001

fluorescent probe made of 24 loops. It is $\tau_{min}^{\text{buff}} = 72$ s and it prevents direct observation of the activation events [19].

To understand the details of the regulatory process that controls mRNA expression we need to quantify the statistics of the activation and inactivation times, as has been performed in cell cultures [11, 12, 14, 15]. However the very short duration of the cell cycles (6–15 minutes for cell cycles 11-13) in early fly development prevents accumulation of statistics about the inactivation events and interpretation of these distributions. Direct observation of the traces suggests that transcription regulation is not static but displays bursts of activity and inactivity. However the eye can often be misleading when interpreting stochastic traces. In this paper we develop a statistical analysis of time dependent gene expression traces based on specially

designed autocorrelation functions to investigate the dynamics of transcription regulation. This method overcomes the analysis difficulties resulting from naturally short traces caused by the limited duration of the cell cycles that make it impossible to infer the properties of regulation directly from sampling the activation and inactivation time statistics. Combining our analysis technique with models of transcription initiation, estimates of the precision of the transcriptional readout and high resolution microscopy imaging of the MS2 cassette under the control of the hunchback promoter in heterozygous flies, we find evidence suggesting bursty transcription initiation in cell cycles 12-13. For the switching timescales we observe experimentally, the autocorrelation function analysis alone is not able to reliably distinguish between different models for promoter activation and we use information about the precision of the transcriptional readout to conclude that transcription is most likely bursty. Based on the analysis of the time traces, we show that the precision of the transcriptional readout in each cell cycle is relatively imprecise compared to the expected precision of the mRNA measurement obtained from fixed samples, both in terms of cell-to-cell variability [5] and embryo-toembryo variability [6]. We discuss the limitations of the inference for models of different complexity in different parameter regimes.

Results

Characterizing the time traces

Before we present our results, we first analyze the traces and present a new analysis technique. We study the transcriptional dynamics of the hunchback promoter (depicted in Fig 1A and 1B) by generating embryos that express an MS2 reporter cassette under the control of the proximal hunchback promoter (Fig 1C), using previously developed tools [17, 18], with an improved MS2 cassette [20] (see Materials and Methods for details). The MS2 cassette was placed towards the 3' end of the transcribed sequence and contained 24 MS2 loop motifs. While the gene is being transcribed, each newly synthesized MS2 loop binds MCP-GFP (expressed at low levels and freely diffusing in the embryo). In each nucleus, where transcription at this reporter is ongoing, we observe a unique bright fluorescent spot, which corresponds to the accumulation of several MCP-GFP molecules at the locus (Fig 1C). We assume that the fluorescent signal from a labeled mRNA disappears from the recording spot when the RNAP reaches the end of the transgene. With this setup we image the total signal in four fly embryos using confocal microscopy, simultaneously in all nuclei (Fig 1D) from the beginning of cell cycle (cc) 11 to the end of cell cycle 13. In each nucleus we obtain a signal that corresponds to the temporal dependence of the fluorescence intensity of the transcriptional process, which we refer to as the time trace of each spot. A cartoon representation of such a trace resulting from the polymerase activity (Fig 1E) dictated by the promoter dynamics (Fig 1B) is shown in Fig 1F. We present examples of the traces analyzed in this paper in S1 Fig and the signal preprocessing steps in the Materials and Methods and SI Section A.

To characterize the dynamics of the *hunchback* promoter we need to describe its switching rates between ON states, when the gene is transcribed by the polymerase at an enhanced rate and the OFF states when the gene is effectively silent with only a small basal transcriptional activity (Fig 1A and 1B). Estimating the ON and OFF rates directly from the traces is problematic due to the buffering time and to the high background fluorescence levels coming from the unbound MCP-GFP proteins that make it difficult to distinguish real OFF events from noise. To overcome this problem, we consider the autocorrelation function of the signal. To avoid biases from different signal strengths from each nucleus, we first subtract the mean of the fluorescence in each nucleus, $F(t_i) - \langle F(t_i) \rangle$ and then calculate the steady state connected autocorrelation function of the fluorescence signal (equivalent to a normalized auto-covariance), $C(\tau)$,

at two time points separated by a delay time τ , $F(t_i)$ and $F(t_i + \tau)$, normalized by the variance of the signal over the traces, according to Eqs 12 and 13 in Materials and Methods. We limit our analysis to the constant expression part of the interphase (which we call "steady state"—we discuss this assumption at the end of the Simulated data Results section) by taking a window in the middle of the trace to avoid the initial activation and final deactivation of the gene between the cell cycles (see <u>Materials and Methods</u> and S2 Fig). We will always work with the *connected* autocorrelation function, which indicates that the mean of the signal is subtracted from the trace. The autocorrelation function is a powerful approach since it averages out all temporally uncorrelated noise, such as camera shot noise or the instantaneous fluctuations of the fluorescent probe concentrations.

Fig 2A compares the normalized connected autocorrelation functions calculated for the steady state expression in the anterior of the embryo (excluding the initial activation and final deactivation times after and before mitosis) in cell cycles 12 and 13 of varying durations: ~ 3 and \sim 6 minutes. Fig 2B shows the same functions for traces that have been curtailed to all have equal length. The steady state signal from cell cycle 11 did not have enough time points to gather sufficient statistics to calculate the autocorrelation function. As expected, the functions decay showing a characteristic correlation time, then reach a valley at negative values before increasing again. Since the number of data points separated by large intervals is small the uncertainty increases with τ . Autocorrelation functions calculated for very long time traces have neither the negative valley nor the increase at large τ . For example, the long-time connected autocorrelation functions calculated from the simulated traces (Fig 2C) of the process described in Fig 1 that are shown in Fig 2D, differ from the short time connected autocorrelation function in Fig 2E calculated from the same trace (see SI Section G for a description of the simulations). As the traces get longer the connected autocorrelation function approaches the longtime results (S4 Fig). The connected autocorrelation function of a finite duration trace of a simple correlated brownian motion (an Ornstein-Uhlenbeck process) displays the same properties (see S5 Fig). The dip is thus an artifact of the finite size of the trace. We also see that the autocorrelation functions shift to the left for short cell cycles (Fig 2A), resulting, for earlier cell cycles, in shorter directly read-off correlation times, defined as the value of τ at which the autocorrelation function decays by e. However, calculating the autocorrelation functions for time traces of equal lengths for all cell cycles (Fig 2B) shows that the shift was also a bias of the finite trace lengths, and after taking it into account, the transcription process in all the cell cycles has the same dynamics (although we note that the dynamics directly read out from this truncated trace is not the true long time dynamics).

This preliminary analysis shows that to extract information about the dynamics of transcription initiation we will need to account for the finite time traces. Additionally, a direct readout of even effective rates from the correlation time is difficult, because the autocorrelation coming from the underlying gene regulatory signal (Fig 1B) is obscured by the autocorrelation due to the timescale needed for the transcription of the sequence containing the MS2 cassette (Fig 1E)—the gene buffering time, τ_{buff} . The observed time traces are a convolution of these inputs (Fig 1F). The analysis is thus limited by the buffering time of the signal ($\tau_{buff} = 72s$ in our system), given as the length of the transcribed genomic sequence that carries the fluorescing MS2 loops divided by the polymerase velocity. A direct readout of the switching rates is only possible if the autocorrelation time of the promoter is larger than the buffering time.

The form of the autocorrelation function and our ability to distinguish signal from noise also depends on the precise positioning and length of the fluorescent gene [19]. A construct with the MS2 transgene placed at the 3' end of the gene (Fig 3A) gives a differentiable readout of the promoter activity even for two sets of fast switching rates between the active and inactive states. However, in this case the weak signal is hard to distinguish from background





Fig 2. Autocorrelation analysis of fluorescent traces from cell cycles 12-13. (A) Autocorrelation functions for traces of different length caused by the variable duration of the cell cycle. Each autocorrelation function is calculated from one embryo and one cell cycle from traces in the anterior region of the embryos. Reading off the autocorrelation time as the time at which the autocorrelation function decays by a value of *e* would give different values for each trace. The analysis is restricted to the steady state part of the traces (as defined in the text and S2 Fig). The durations of the steady state windows are given in SITable I. (B) Autocorrelation functions calculated for the same traces reduced to having equal trace lengths, all equal to the trace length of the shortest trace (101*s*), show that the differences observed in panel A are due to finite size effects. In the curtailed traces all sequential time points until the 101*s* time point were used. (C) An example of a signal simulated using the process described in Fig 1 for 300 seconds (blue curve) for a two state model. Taking the whole 300 second interval (red dashed lines) gives a good approximation of the average signal (red line) and the effect of finite size on the autocorrelation is strong (E). The sampling rates of the four embryos are: 13.1*s*, 10.2*s*, 5.1*s* and 4.3*s*, respectively. Parameters for the simulation in (C-E) are: $k_{on} = k_{off} = 0.06s^{-1}$, sampling time dt = 4s, for the red curve T = 300s and M = 2000 nuclei, for the green curve T = 60s and M = 10000 nuclei (same total amount of data). These parameters were chosen for illustrative purposes.

doi:10.1371/journal.pcbi.1005256.g002

fluorescence levels. Conversely, a 5' positioning of the transgene (Fig 3B) is insensitive to background fluorescence. However it only differentiates autocorrelation functions calculated from very slow switching processes [19]. In summary, a construct with the MS2 placed at the 3' end of the gene allows for a direct readout of the transcriptional kinetics in a much wider range of switching rates than a 5' construct, although the autocorrelation function of a 3' construct is more sensitive to background fluorescence.

Method development—Promoter switching models

The promoter activity we are interested in inferring can in principle be described by models of varying complexity (see Fig 1A). We consider and compare three types of models in this

COMPUTATIONAL BIOLOGY

hb aen

Two state simulation

Three state simulation

Three state prediction

- Two state prediction

PLOS

Α

С

1



hb gene

3' simulation

3' prediction

5' simulation

-5' prediction

Pol II



8 999



В

D

5

1

doi:10.1371/journal.pcbi.1005256.g003

paper. We note this is a small subset of possible models. In particular, we do not consider models with multiple levels of transcription as was considered in [21] or reversible promoter cycles. In the simplest case, the gene is consecutively yet noisily expressed. The RNAP starts transcribing following a Poisson distribution of discrete ON-activation (or firing) events-this has previously been called a static promoter (not represented in Fig 1A). After the polymerase binds, the next polymerase cannot bind before the promoter is cleared (a timescale estimated to be $\tau_{\rm block} \sim 6s$ in our experiments). The effective firing rate of this model is the Poisson rate, r, shifted by a deterministic $\tau_{block} \sim 6s$, $r_{eff} = (\tau_{block} + r^{-1})^{-1}$, and we call this discrete time model a Poisson-like promoter. Although the promoter dynamics would be uncorrelated in this case, the gene buffering would still produce a finite correlation time (see SI Section F). Alternatively, the promoter could have two well defined expression states: an ON state during which the polymerase is transcribing at an enhanced level and an OFF state when it transcribes at a basal level. This situation can be modeled by stochastic switching between the two states with rates $k_{\rm on}$ and $k_{\rm off}$ (left panel in Fig 1A and Materials and Methods). However, as was previously observed in both eukaryotic and prokaryotic cell cultures [11, 12, 14, 15], once the gene is switched off the system may have to progress through a series of OFF states before the gene can be reactivated. Recently these kinds of cycle models have been discussed for the hunchback promoter [22]. The intermediate states can correspond to, for example, the assembly of the transcription initiation complex, opening of the chromatin or transcription factor cooperativity. These kinds of situations can either be modeled by a promoter cycle (middle panel in Fig 1A and Materials and Methods), with a number of consecutive OFF states, or by an effective two state model that accounts for the resulting non-exponential, but gamma function distribution of waiting times in the OFF state (right panel in Fig 1A and Materials and Methods). The time the polymerase spends transcribing the DNA does not dependent on the promoter model.

In both the two-state and promoter cycle model the gene switches from the ON to the OFF state with exponentially distributed waiting times described by a rate k_{off} (Fig 1A). In the two-state model the jumps from the OFF to the ON state are also exponentially distributed with a switching rate k_{on} (Fig 1A). In the three state cycle model considered in this paper, an inactive gene can be in two different OFF states. The gene leaves these states with different switching rates, k_1 and k_2 , respectively. The ordering of k_1 and k_2 is impossible to detect in the current experiment. In the three state cycle model we can define an effective on-switching rate $k_{on}^{eff} = (1/k_1 + 1/k_2)^{-1}$. k_{on}^{eff} corresponds to the inverse of the average waiting time in the overall OFF state, and the waiting times for exiting this effective OFF state are not exponentially distributed. The gamma function distributed switching time is an approximation of this effective rate. We present our method for all of these models and consider all but the gamma function distributed switching time is of *hunchback* promoter dynamics.

Method development—Autocorrelation approach

To infer the transcription dynamics from the data we built a mathematical model that calculates the autocorrelation functions accounting for the experimental details of the probes, incorporating the MS2 loops at various positions along the gene and correcting for the finite length of the time traces. The basic idea behind our approach is that while the initiation of transcription is stochastic and involves switching between the ON and possibly a number of OFF states (X(t) in Fig 1B denotes the binary gene expression state), if we assume a constant elongation velocity the obscuring of the signal by the probe design is completely deterministic [18, 23], which results in the random variable $a(i, t) \in \{0, 1\}$ that describes the presence or absence of the polymerase at position *i* at time *t* (Fig 1E). We count the progression of the polymerase in discrete time steps, where one time step corresponds to the time it takes the polymerase to cover a distance of 150 base pairs equal to its own length (Fig 4A). The promoter dynamics can thus be learned from the noisy autocorrelation function of the fluorescence intensity normalized by the intensity coming from one MS2 loop, $F(t) = \sum_{i=1}^{r} L_i a(i, t)$ (Fig 1F), even for switching timescales smaller than the fluorescent probe buffering time τ_{buffb} provided the parameters of the probe design encoded in the loop function L_i (positioning of the probe etc.) are known (Fig 1C) and the intensity signal is calibrated knowing the fluorescence intensity coming from one MS2 loop [18].

Broadly, our model assumes that once the promoter is in an ON state the polymerase binds and deterministically travels along the gene producing MS2 loops containing mRNA that immediately bind MCP and result in a strong localized fluorescence (Fig 4). The presence or absence of a polymerase at position i at time t, a(i, t) is simply a delayed readout of the promoter state at time t - i, a(i, t) = X(t - i) where t is measured in polymerase time steps (Fig 1B). We assume that polymerase is abundant and that at every time step a new polymerase starts transcribing, provided the gene is in the ON state (Fig 1B and 1E). The amount of fluorescence produced by the gene at one time point is determined by the number of polymerases on the gene (Fig 4A). The amount of fluorescence from one polymerase that is at position i on the gene depends on the cumulated number of loops that the polymerase has produced L_{ij} where $1 \le i \le r$, *r* corresponds to the *maximum* number of polymerases that can transcribe the gene at a given time and $L_i = 1$ corresponds to one loop fluorescing, as depicted in the cartoon in Fig 4B. The known loop function L_i depends on the build and the position of the MS2 cassette on the gene, it is input to the model and does not necessarily take an integer value since the polymerase length and the loop length do not coincide (Fig 4B). Given the average fraction of time the transcription initiation site is occupied by the polymerase, $P_{\rm on}$, the average fluorescence in the steady state is:

$$\langle F \rangle = P_{\rm on} \sum_{i=1}^{r} L_i. \tag{1}$$

Since we assume the polymerase moves deterministically along the gene, seeing a fluorescence signal both at time *t* and position *i*, and at time *s* and position *j* means the gene was ON at time t - i and s - j, which is determined by how many loops (*i* and *j*) the polymerase has produced. Taking the earlier of these times, we need to calculate the probability that the gene is also ON at the later time. The autocorrelation function of the fluorescence can thus be written as:

$$\langle F(t)F(s)\rangle = \sum_{i=1}^{r} \sum_{j=1}^{r} L_i L_j P(\text{gene was ON at time})$$
$$\min(t-i,s-j)) \cdot A(|t-i-s+j|), \tag{2}$$

where A(n) is the probability that the gene is ON at time *n* given that it was ON at time 0, and time is expressed in polymerase steps. The precise form of P_{on} , P(gene was ON at time min(t - i, s - j)) and A(|t - i - s + j|) depends on the type of the promoter switching model. We assume that the polymerase moves at constant speed along the gene and that there is no splicing throughout the transcription process. We give explicit expressions for all the models used in the Materials and Methods section and the Supplementary Information. Knowing the design of the construct (length of the probe and number of loops that have been transcribed at each position) and calibrating the signal, we use Eq 1 to directly learn P_{on} from the data. In the two and multistate models P_{on} provides us with the ratio of switching rates and we then use Eq 2 to obtain their particular values (see Materials and Methods).





Fig 4. The gene expression model used in the autocorrelation function calculation. The autocorrelation inference approach is based on the idea that the stochastic transcriptional dynamics can be deconvoluted from the signal coming from the deterministic fluorescent construct, if we know the gene construct design. (A) A concatenation of snapshots of the gene from *r* consecutive time steps. A polymerase covers a length on the gene corresponding to its own length in one time step, producing about two MS2 loops. The gene has total length *r* and at any position *i* along the gene $L_i < 24$ loops have been produced. (Top right) The promoter state as a function of time and (center right) an instantaneous snapshot of the gene corresponding to transcription from this promoter. (B) The construct design is encoded in the loop function L_i . As the polymerase moves along the gene it produces MS2 loops. L_i is an average representation in terms of polymerase time steps of how many loops have been produced by a single polymerase. It is based on the experimental design shown on the left of the panel.

doi:10.1371/journal.pcbi.1005256.g004

To avoid biases coming from nucleus to nucleus variability, we calculated the normalized connected correlation function defined in Eqs 12 and 13 in Materials and Methods. The theoretically calculated connected autocorrelation function, C_r (Eq 14, which corresponds to the longtime correlation function in Fig 2C and 2D), differs from the empirically calculated connected autocorrelation from the traces, c(r) (Eqs 12 and 13 in Materials and Methods, which correspond to the short time correlation function in Fig 2C and 2E), due to finite size effects coming from spurious correlations between the empirical mean and the data points. Since by definition the mean of a connected autocorrelation function is zero (see Eqs 12 and 13 in Materials and Methods), the area under the autocorrelation function must be zero. For short traces this produces the artificial dip discussed in Fig 2, which for long traces is not visible, as it is equally distributed over long times. To compare our theoretical and empirical correlation functions we explicitly calculate the finite size correction and include this correction in our analysis (Materials and Methods and SI Section H and I).

In this paper, we have analyzed data from fly embryos with 3' promoter constructs only, limiting ourselves to the steady state part of the trace (see S2 Fig). However the method can also be applied to non-steady state systems (see SI Section C) and to other constructs, including cross-correlation functions calculated from signals of different colors inserted at different positions along the gene (see SI Section J). We use simulated data to show that prediction and inference are possible for cross-correlation functions of a two-colored signal (see S6 Fig), but that the accuracy of inference is limited by the use of the 5' probe.

Simulated data

To check that the inference method correctly infers the parameters of the model, we first tested the autocorrelation based inference on simulated short-trace data with underlying molecular models with different levels of complexity (Fig 3C) for a construct with the MS2 probe placed at the 3' end of the gene (Fig 3A). In Fig 3D we compare autocorrelation functions for the three state model for constructs with the MS2 loops positioned at the beginning of the transcribed region (5', Fig 3B) and at the end of the transcribed region (3', Fig 3A), and the cross-correlation function calculated from a two-colored probe construct (Fig 3E). The analytical model correctly calculates the short trace autocorrelation function and is able to infer the dynamics of promoter switching for all models. It can also be adapted to infer the promoter switching parameters for any intermediate MS2 construct position, given the limitations of each of the constructs discussed above [19].

The autocorrelation function based inference reproduces the underlying parameters of the dynamics with great accuracy not just for switching timescales longer than the gene buffering time, τ_{buff} , that obscures the signal (Fig 3F), but also for smaller timescales that are within an order of magnitude of the gene buffering time. In Fig 3F we show the results of the inference for the 3' two state model for different values of the ON and OFF rates, k_{on} and k_{off} . For switching timescales much shorter than the gene buffering time, the autocorrelation function coming from the length of the construct dominates the signal and the precision of the inference goes down. For very fast switching rates ($k_{on} + k_{off} > 0.12s^{-1}$), increasing the length of the traces or the number of nuclei (red vs blue curve for values of $k_{on} + k_{off}$ larger than $0.1s^{-1}$ in Fig 3F) does not help estimate the properties of transcription. In this regime, the inferred value of $k_{\rm on} + k_{\rm off}$ disagrees with the true parameters even when the inference uses long time traces and a large number of nuclei. For intermediate switching rates $(0.07 - 0.12s^{-1})$, increasing the trace length or increasing the number of nuclei extends the inference range (black and green dashed lines vs blue solid line in Fig 3F), and in all cases increasing the number of nuclei decreases the uncertainty as can be seen from the smaller error bars (shown only for the red and blue lines for figure clarity).

Using two colored probes attached at different positions along the gene gives two measurements of transcription and allows for an independent measurement of the speed of the polymerase—one of the parameters of the model that currently must be taken from other experiments. While the estimates of polymerase speed in the fly embryo are reliable [18], this parameter has been pointed out as a confounding factor in other correlation analyses [24].

The autocorrelation approach also correctly infers the parameters of transcriptional processes when applied to traces that are out of steady state (see SI Section C). However, since the process is no longer translationally invariant more traces are needed to accumulate sufficient statistics. For this reason, in the current analysis of fly embryos we do not analyze the transient dynamics at the beginning and end of each cycle and we restrict ourselves to the middle of the interphase assuming steady state is reached (see <u>S2 Fig</u> for details). We do not know whether the underlying dynamics is completely in steady state. We limit our analysis to a time frame window where the intensity of the fluorescence signal plateaus (see S2 Fig for an example). We can motivate the steady state assumption *a posteriori*: the inferred switching timescales (smaller than 50s) are small enough for the system to relax to steady state within one cell-cycle. However we cannot fully rule out other mechanisms that could keep the system out of steady state (such as changes in the Bicoid concentration).

Fly trace data analysis

We divided the embryo into the anterior region, defined as the region between 0% and 35% of the egg length (the position at 50% of the egg length marks the embryo midpoint), where *hunchback* expression is high, and the boundary region, defined as the region between 45% and 55% egg length, where *hunchback* expression decreases. The mean probability for the gene to be ON during a given cell cycle P_{on} (restricted to the times excluding the initial activation and deactivation of the gene, which we will call the steady state regime), given by Eq 1, is consistent between the four embryos in cell cycle 12 and 13, both in the anterior region and at the boundary (Fig 5A). The probability for the gene to be ON is over three fold higher in the anterior region than in the boundary and does not change with the cell cycle. $P_{on} \sim 0.5$ in the anterior indicates that in each nucleus the polymerase spends about half the steady state expression time transcribing the observed gene. At the boundary the gene is transcribed on average during about 10% of the steady state part of the cell cycle. The estimates for P_{on} in the earlier cell cycles were not reproducible between the four embryos, likely because the time traces were too short to gather sufficient statistics to accurately calculate the maximum and average of the signal. We concentrated on cell cycle 12 and 13 for the remainder of the analysis.

Initial comparison of the promoter models. Based on the different behavior at the boundary and in the anterior, we separately inferred the transcriptional dynamics parameters in the two regimes, using the autocorrelation approach that corrects for finite time traces for different models. The Poisson-like promoter model, the two and three state cycle models all provide reasonably good fits to all the traces in both regions (see Fig 5B for an example and S3 Fig for the fits in both regions in all embryos). Although the fit of the Poisson-like promoter model (red line) only captures the short time behavior of the measured autocorrelation function, there is not enough statistical evidence from the autocorrelation analysis to exclude this model. The two and three state model fits are indistinguishable (Fig 5B) and the two state fit is reproducible between cell cycles and embryos (Fig 5C and 5D). The form of the autocorrelation function in the Poisson-like promoter model is completely determined by the autocorrelation signal of the fluorescent construct (the loop function L_i), since the random firing is itself uncorrelated. In the two and three state models, the autocorrelation signal from the fluorescent construct is convoluted with the autocorrelation signal of the promoter. The fact that the Poisson-like promoter model fits the data so well, indicates that the autocorrelation time of the promoter is comparable to the autocorrelation time of the fluorescent construct. In S7 Fig we plot the autocorrelation functions for simulated two state models with different correlation times $(k_{on} + k_{off})$ and a Poisson-like promoter with the same P_{on} . For short correlation times, the Poisson-like promoter model and two state model have indistinguishable autocorrelation functions, just like in the analyzed data. For long autocorrelation times, the difference between the two models is clear.

The three state fit is reproducible at the level of the sum of the effective ON and OFF rates (same fit as shown for the two state model in Fig 5C), but gives fluctuating values for k_1/k_2 , the parameter ratio determining how well it is approximated by a two state model (see S8 Fig, $k_1/k_2 < 1$ describes one fast reaction between the OFF states, effectively giving a two state model, while $k_1/k_2 = 1$ gives equal weights to the two reactions, clearly distinguishing two OFF





doi:10.1371/journal.pcbi.1005256.g005

COMPUTATIONAL BIOLOGY

PLOS

states). Since the two state model is reproducible, and has lesser complexity we will further consider only the two state and Poisson-like promoter models.

Discussion of the two state model. For the two state model, the inference procedure independently fits the characteristic timescale of the process from the autocorrelation function, defined as the inverse of the sum of two rates, $k_{on} + k_{off}$ (Fig 5C), and then uses an independent

fit of the probability of the gene to be ON, $P_{\rm on} = k_{\rm on}/(k_{\rm on} + k_{\rm off})$, to disentangle the two rates (Fig 5D). Examples of the simulated promoter state over time with the rates' inferred values are shown in Fig 5E (for the anterior region) and Fig 5F (for the boundary region). Assuming the two state model we find that the characteristic timescale, $(k_{on} + k_{off})^{-1}$, in most embryos is slightly shorter at the boundary ($\sim 25s$) than in the anterior region ($\sim 33s$) and the variability between the two cell cycles is comparable to the embryo-to-embryo variability (Fig 5C). Both timescales are much larger than the polymerase blocking time, $\tau_{block} \sim 6s$, during which a second polymerase cannot bind because the first one has not cleared the binding site (shown as the gray dashed line in Fig 5D), which sets a natural scale for the timescales we can infer. We find that in the anterior region of the embryo the two switching rates, k_{on} and k_{off} , show variability from embryo to embryo (between $0.009s^{-1}$ to $0.078s^{-1}$ —see Table I and II in the SI), but always scale together, which gives the observed one-half probability of the gene to be ON in a given nuclei during the steady state part of the interphase. Since the polymerase in the anterior on average spends half the steady state interphase window transcribing the gene, the inferred rates suggest a clear bursting behavior of the transcription process, with switching between an identifiable active and inactive state of the promoter if the two state promoter is correct, and rare firing if the Poisson-like promoter model is correct.

At the boundary k_{on} is much smaller than in the anterior with very little embryo to embryo variability, while k_{off} has a similar range in the anterior and at the boundary. This behavior is expected since high Bicoid concentrations in the anterior upregulate the transgene whereas lower concentrations at the boundary result in smaller activation rates. The ratio of the average k_{on} rates at the boundary and anterior is ~ 5, which can be compared to the 4 fold decrease expected from pure Bicoid activation, assuming the Bicoid gradient decays with a length scale of $100\mu m$ [25] and comparing the activation probabilities in the middle of the anterior and boundary regions. Given the crudeness of the argument stemming from the variability of the Bicoid gradient in the boundary region and the uncertainty of the inferred rates, these ratios are in good agreement and suggest that a big part of the difference in the transcriptional process between the anterior and boundary is due to the change in Bicoid concentration. Of course other factors, such as maternal Hunchback, could also affect the promoter, leading to discrepancies between the two estimates.

Discriminating between two and three state models. The current data coming from four embryos and ~ 50 nuclei in each region with trace lengths of $\sim 300s$ does not make it possible to distinguish between the two and three state models. We asked whether having longer traces or more nuclei could help us better characterize the bursty properties. We performed simulations with characteristic times similar to those inferred from the data $(k_{on}^{eff} + k_{off} = 0.01)$ assuming a two state $(k_{on}^{eff} = k_{on}, \underline{Fig 6A})$ and three state cycle model (Fig <u>6B</u>). We then inferred the sum of the ON and OFF rates $(k_{on}^{eff} + k_{off})$ and the ratio of the two ON rates (k_1/k_2) . If the two ON rates are similar $(k_1/k_2 \sim 1)$, we infer a three state model. If one of the rates is much faster $(k_1/k_2 \sim 0)$, which implies $k_{on}^{eff} = k_{on}$, we infer a two state model. We find that having more nuclei, which corresponds to collecting more embryos, would not significantly help our inference. However looking at longer traces would allow us to disambiguate the two scenarios, if the traces were 4 times longer, or ~ 20 minutes long. Since cell cycle 14 lasts \sim 45 minutes, analyzing these traces could inform us about the effective structure of the OFF states. However in cell cycle 14, several other direct Bicoid targets (mostly transcription factors) are likely to be expressed, so additional regulatory elements could be responsible for the observed transcriptional dynamics compared to cell cycle 12 and 13. Our results suggest that with our current trace length we should be able to identify a two state model with large certainty, but we could not clearly identify a three state



Fig 6. Longer time traces help distinguish between two state and three state cycle models. A. Inference from data generated by a two state model, which corresponds to $k_1/k_2 = 0$, from traces of different lengths *T* and using different numbers of nuclei *N* shows that longer traces help increase the probability to correctly learn the model type. Increasing the number of nuclei for short traces shows little improvement. The inference is repeated 50 times per condition. The experimental conditions studied in this paper are closest to the T = 240s and N = 50 nuclei panel. B. The same numerical experiment but assuming a three state cycle model, which corresponds to $k_1/k_2 = 1$. Parameters of the simulations: $P_{on} = 0.1$, $k_{off} + 1/(1/k_1 + 1/k_2) = 0.02s^{-1}$ and $k_1/k_2 = 0$ in A, and $k_1/k_2 = 1$ in B.

doi:10.1371/journal.pcbi.1005256.g006

model. Our data may point towards a more complex model than two state, but different kinds of multistate models or a two state model obscured by other biases cannot be ruled out.

The error bars for the autocorrelation functions (see Fig 5B and S3 Fig) describe the variability between nuclei coming from both natural variability and measurement imprecision. While the autocorrelation function is insensitive to white noise, it does depend on correlated noise. The noise increases for large time differences τ , as the number of pairs of time points that can be used to calculate the autocorrelation function decreases and in our inference we reweigh the points according to their sampling, so that the noise does not impair the precision of our inference. The error bars on the inferred parameters (e.g. Fig 5A, 5C and 5D) are due to variability between nuclei and are obtained from sampling different subsets of the data in each region and cell cycle. Additionally to the inter-nuclei and experimental noise, there is natural variability between embryos. Since each nucleus transcribes independently and we assume similar Bicoid concentrations in each of the regions, the inter-embryo variability is of a similar scale as the inter-nuclei variability (Fig 5C), as one expects given that the Bicoid gradient is reproducible between embryos [25]. Additonally, variability between Bicoid gradients in different embryos, for example due to their different lengths [26], could also contribute to the observed variability.

Accuracy of the transcriptional process

At the boundary, neighboring nuclei have dramatically different expression levels of the Hunchback protein. From measurements of the Bicoid gradient, Gregor and collaborators estimated that for two neighboring nuclei to make different readouts, they must be able to distinguish Bicoid concentrations that differ by 10% [27]. Following the Berg and Purcell [28]

argument for receptor accuracy, and using measurements of diffusion constants for Bicoid proteins from cell cycle 14, the authors showed that, based on protein concentrations, the *hunchback* gene is not able to read-out the differences in the concentrations of Bicoid proteins to the required 10% accuracy in the time that cell cycle 14 lasts. Even considering the revised higher values of Bicoid's diffusion coefficient measured in a subsequent study [5], the precision of the Bicoid gradient read-out remains difficult to explain. The authors invoked spatial averaging of Hunchback proteins as a possible mechanism that achieves this precision. Spatial averaging can increase precision, but it can also smear the boundary. Erdmann et al calculated the optimal diffusion constant Hunchback proteins must have for the averaging argument to work [29] and showed it is consistent with experimental observations [5, 25]. However precision can already be established at the mRNA level and, using measurements on fixed embryos, Little and co-workers found that the relative intrinsic nuclei-to-nuclei variability of the mRNA reduced this variability to ~ 10% [6].

Here we go one step further and use our direct measurements of transcription from the *hunchback* gene to directly estimate the precision with which the *hunchback* promoter makes a readout of its regulatory environment in a given cell cycle in a given region of the embryo, $\delta P_{\rm on}/P_{\rm on}$. $\delta P_{\rm on}/P_{\rm on}$ is the relative error of the probability of the gene to be ON averaged over the steady state part of a cell cycle. Since the total number of mRNA molecules produced in a given cycle is proportional to $P_{\rm on}$ (shown in S9A Fig as a function of embryo length), the precision at the level of *produced* mRNA in a given cycle is equal to the precision in the expression of the gene, $\delta mRNA/<mRNA > = \delta P_{\rm on}/P_{\rm on}$. The accuracy of transcription activation is encoded in the stochasticity of gene activation.

In the two state model, the gene randomly switches between two states: active and inactive, making a measurement about the regulatory factors in its environment and indirectly inferring the position of its nucleus. Since no additional information is provided by a measurement that is strongly correlated to the previous one, the cell can only base its positional readout on a series of independent measurements. Two measurements are statistically independent, if they are separated by at least the expectation value of the time τ_i it takes the system to reset itself:

τ

$$_{i} \sim \frac{1}{k_{\rm on}^{\rm eff} + k_{\rm off}},$$
(3)

where in a two state model $k_{on}^{eff} = k_{on}$. A more detailed estimate obtained by computing the variance of the time spent ON by the gene during the interphase (see SI Section K) shows that Eq 3 underestimates the time needed to perform independent measurements. We find that for a two state model the accuracy of the readout of the total mRNA produced is limited by the variability of a two state variable divided by the estimated number of independent measurements within one cell cycle:

$$\frac{\delta \text{mRNA}}{\langle \text{mRNA} \rangle} = \sqrt{2 \frac{\tau_i (1 - P_{\text{on}})}{T P_{\text{on}}}},\tag{4}$$

where *T* is the duration of the cell cycle and the factor $\sqrt{2}$ is a prefactor correction to the naive estimate. Eq 4 is valid in the limit of $T >> \tau_i$ (the exact result if given in SI Section K). Using the rates inferred from the autocorrelation analysis (Fig 5D) we see that the precision of the gene readout is much lower at the boundary than in the anterior, does not change with the cell cycle and is reproducible between embryos (blue points on the ordinate in Fig 7A). In the anterior part of the embryo it reaches ~ 50%, while at the boundary, it is very large, ~ 150%, even at cell cycle 13.



Fig 7. Precision of the *hunchback* **gene transcription readout.** A. Comparison of the relative error in the mRNA produced during the steady state of the interphase estimated empirically from data (abscissa) and from theoretical arguments in Eq.4 for a two state switching promoter (blue symbols) using the inferred parameters in Fig 5C (ordinate), and theoretical arguments in Eq.5 for a Poisson-like static promoter (red symbols) using the inferred parameters in Fig 5A (ordinate), in the anterior (circles and squares) and the boundary (diamonds and triangles) regions. The theoretical prediction for the two state promoter shows very good agreement with the data, whereas the Poisson-like promoter shows poor agreement, especially in the boundary region. B. The relative error in the total mRNA produced in cell cycle 13 directly estimated from the data as the variance over the mean of the steady state mRNA production (red line, same data as in A), sum of the intensity over the whole duration of the interphase (blue line) and the total mRNA produced during cell cycles 11 to 13 (green line) for equal width bins equal to 10% embryo length at different positions along the AP axis. Each line describes an average over four embryos (see S9C Fig for the same data plotted separately for each embryo) and the error bars describe the variance. To calculate the total mRNA produced ouring the steady state for a two state, $k_1/k_2 = 0$, (solid lines) and thrace back their lineage through cycle 12 to cycle 11. We then sum the total intensity of each nuclei in cell cycle 13 and half the total intensity of its mother and 1/4 of its grandmother. C. Comparison of the relative error rate for the Poisson-like and three state cycle model, $k_1/k_2 = 1$, (dashed lines) with the same $k_{on}^{eff} + k_{off}$, for different values of k_{on}^{eff} and k_{off} , shows that the three state cycles system allows for greater readout precision. D. A comparison of the theoretical prediction of the steady state relative error rate for the Pois

doi:10.1371/journal.pcbi.1005256.g007

In the Poisson-like promoter model, to calculate the relative error in the total mRNA produced, the polymerase arrival times are described by an effective firing rate of $r^{\text{eff}} = (\tau_{\text{block}} + 1/r)^{-1}$. Within this model, the fraction of the total time the polymerase cannot bind, because the binding site is occupied is $P_{\text{on}}^{\text{p}} = \tau_{\text{block}} \cdot r^{\text{eff}}$ (see SI Section F). The total produced mRNA is then proportional to the time the gene is transcribed, and the relative error in the total produced mRNA depends on the relative error of the firing times of this modified Poisson process and the number of independent measurements, $n_P = T/(\tau_{\text{block}} + 1/r)$ (see SI Section K):

$$\frac{\delta \text{mRNA}}{\langle \text{mRNA} \rangle} = \sqrt{\frac{\tau_{\text{block}}(1 - P_{\text{on}}^{\text{p}})}{T}}.$$
(5)

Using the rates for P_{on}^{p} inferred from the data (Fig 5A), the relative error in the total mRNA produced (blue points on the ordinate in Fig 7A) is slightly higher in the boundary region (~ 15%) than in the anterior of the embryo (~ 10%) and does not change with the cell cycle.

We can compare both of these theoretical estimates with direct estimates of the relative error of the total mRNA produced during a cell cycle, δ mRNA/mRNA, from the data. We divide the embryo into anterior and boundary strips, as we did for the inference procedure and calculate the mean and variance of $P_{\rm on}$. These empirical estimates of the gene measurement precision agree with the theoretical estimates (Fig 7A) for the two state model, but disagree with the predictions of the Poisson-like promoter model. For completeness, we also calculated the relative error for the three state model (see S9B Fig), which shows better agreement than the Poisson-like promoter model but slightly worse than a two state model. We verified that our conclusions about the scale of our empirical estimates are the same for all embryos (S9C Fig) and do not depend on the definition of the boundary and anterior regions (S9D Fig). Since the predicted relative error for the Poisson-like promoter model is lower than the relative error calculated directly from the data, we could imagine that the data estimate is more susceptible to additional sources of experimental noise. However the very large disagreement between the Poisson-like promoter prediction and the data at the boundary suggests the Poisson-like promoter model is not an accurate description. This higher experimental variability also cannot be explained by variable expression levels within the regions we are considering: S9D Fig shows that the experimental relative error does not significantly decrease if we take smaller windows and taking the P_{on} in the boundary region to range from 0.35 to 0.5 (Fig 3A) would translate into a, at best, two fold increase in the relative error predicted from the Poisson-like promoter model, which is not enough to reach the experimentally observed relative error. While we are unable to rule out the Poisson-like promoter model based on the fit to the autocorrelation function, a different statistic—the relative error in the produced mRNA suggests that the promoter is most likely well described by a two state model, and possibly a three state cycle.

To see whether temporal integration of the mRNA produced can increase precision, we compared the empirical estimate of the steady state mRNA production (red line in Fig 7B) to the relative error of the total mRNA produced in cell cycle 13 (blue line in Fig 7B) and the total mRNA produced from cell cycle 10 to 13 (green line in Fig 7B) averaged over embryos. Assuming that the mRNA molecules are equally divided between daughter cells during division, and they are all kept in the cell throughout cell-cycles 10–13 (which is incorrect but provides a best case estimate), then each nuclei has the total mRNA produced in cell cycle 13, 1/2 of the total mRNA produced by its mother in cell cycle 12, 1/4 of the mRNA produced by its grandmother in cell cycle etc. While we see about a 1/3 increase in the precision at the

boundary from integrating the mRNA produced in different cell cycles, the estimate in the anterior region is not helped by integration over the cell cycles.

For completeness of the discussion of the relative errors in the different models, we calculated the relative error assuming the same $k_{off} + k_{on}^{eff}$ for a three state cycle $(k_{on}^{eff} = k_{1}^{-1} + k_{2}^{-1})^{-1}$ as for a two state model $(k_{on}^{eff} = k_{on})$ for different values of k_{off} and k_{on}^{eff} (Fig 7C). We found that the relative error is always lower for the three state cycle model and the error decreases regardless of the duration of the cell cycle. As expected from Eq 4, the relative error is decreased by increasing k_{on} and decreasing k_{off} . However the increase in precision from a three state cycle model in the parameter regime we inferred for the two state model in the fly embryo is relatively modest (from $\sim 74\%$ for the two state model to $\sim 67\%$ for the three state model). Similarly, in Fig 7D we compared the prediction for the relative errors for the Poisson-like promoter model to two state models with the same probability of the gene to be transcribed, P_{on} , but different switching rates between the two states (k_{on} and k_{off}). Faster switching increases the precision of the two state promoter, since the number of independent measurements increases. The Poisson-like promoter is always more accurate than the two state promoter.

Many previous analysis of precision from static images calculated the relative error of the distribution of a binary variable, which in each nucleus was 1 if the nucleus expressed mRNA in the snapshot, and 0 if it did not express [30, 31]. We analyzed our data using this definition of activity (see S9E Fig for mean activity as a function of position) and found that for most embryos the relative error in the anterior drops to zero (S9F Fig), indicating that all nuclei in a given region show the same expression state, but at the boundary the precision is still \sim 50%, in agreement with previous reports about the total mRNA in the nucleus [6]. This provides additional evidence for the bursty nature of transcription in the anterior of the embryo, in agreement with previous results that showed a relationship between Bicoid concentration and transcriptional burst of downstream genes [32].

Discussion

In contrast to previous studies [17, 18], including ours, which failed to show evidence for bursty switching of the *hunchback* promoter, by developing more advanced analysis techniques we show that the promoter has distinct periods of enhanced polymerase transcription followed by identifiable periods of basal polymerase activity. Our conclusions are based on combining a new autocorrelation based analysis approach, applied to live imaging MS2 data to infer switching parameters, with an analysis of the precision of readout of promoter. The data we used in this paper was generated with a modified MS2 cassette [20] (see the Experimental procedures section in <u>Materials and Methods</u>) compared to the previously published data [17]. However the difference in our conclusions mainly comes from a detailed analysis of the traces.

Quantification of transcription from time dependent fluorescent traces in prokaryotes and mammalian cell cultures has shown that the promoter states cycle through at least three states [11, 12]. In one of these states the polymerase transcribes at enhanced levels, while in most of the remaining states the transcription machinery gets reassembled or the chromatin remodels. We find that in the living developing fly embryo, the *hunchback* promoter also cycles through at least two states, although based on the parameter inference alone we cannot conclusively rule out the possibility of a static promoter with a Poisson-like firing rate or of a more complex promoter with more effective states when the gene is inactive. Only a combination of the inferred parameters using the autocorrelation function with another statistic (the relative error

in the produced mRNA) allows us to favor the two state model (or more complex models) over the other considered mode of transcription activation.

The main impediment to distinguishing different types of transcriptional models comes from the very short durations of the interphase in the early cell cycles when the *hunchback* gene is expressed. We showed using simulations that increasing the number of embryonic samples would not help us distinguish between two and three state models, however looking at longer time traces would be informative (Fig 6). Since cell cycle 14 lasts about 45 minutes, our analysis shows that the steady state part of the interphase provides enough time to gather statistics that can inform us about the detailed nature of the bursts. Unfortunately, other transcription factors, such as the other gap genes regulating *hunchback* expression in cell cycle 14, could possibly change the nature of the transcriptional dynamics in a time dependent manner. We showed that the transcriptional dynamics is constant and reproducible in the earlier cell cycles (12-13) (Fig 2), so independently of the question of the nature of the bursts, it would be very interesting to see whether and how it changes when the nature of regulation changes.

In the parameter regime of relatively fast switching that we inferred from our data, the autocorrelation function for the two state model and the Poisson-like promoter model are very similar. In this parameter regime, the form of the autocorrelation function is governed by the autocorrelation of the fluorescent probe. So while the autocorrelation function approach is able to disentangle the real promoter switching from the buffering of the construct to determine the parameters assuming an underlying model, we cannot conclusively discriminate between these two models, without looking at other statistics. Using simulations (S7 Fig), we showed that for promoters with slower switching characteristics, this discrimination task is possible and the autocorrelation function approach alone can reliably discriminate between different models. In the parameter regime inferred for the *hunchback* promoter, having longer traces would not be helpful for this discrimination task and we have to look for other statistics (S10 Fig). However using new constructs with MS2 binding sites that have higher binding affinity to MCP and decrease the noise from the binding/unbinding of MCP to the RNA would make it possible to use shorter MS2 cassettes without increasing background fluorescence. These cassettes would decrease the buffering time and extend the parameter regime in which we can distinguish between the Poisson-like and two state promoter models.

Alternatively to focussing on longer traces, a construct with two sets of MS2 loops placed at the two ends of the gene that bind different colored probes could be used to learn more about transcription dynamics [33]. We do not have access to data coming from such a promoter, but our analysis approach can be extended to calculate the cross-correlation function between the intensities of the two colored probes. Such cross-correlation analysis has previously been used to study transcription in cell cultures [34], transcriptional noise [35] and regulation in bacteria [36, 37]. Our theoretical prediction for such a cross-correlation function agrees with simulation results (Fig 3C). Unfortunately, the cross-correlation function with one set of probes inserted at the 5' end and the other at the 3' shares the same problems of a 5' construct. For fast switching rates, such a cross-correlation function suffers from the large buffering time $(\tau_{\text{buff}} \sim 300 \text{ s in [18]})$ drawback of the 5' design and can only be used for inferring large switching rates [38] (see S6 Fig). Similarly, the cross-correlation function cannot discriminate between a two state and Poisson-like promoter for relatively fast switching. However, it does gives us access into dynamical parameters of transcription such as the speed of polymerase and it is able to characterize whether mRNA transcription is in fact deterministic and identify potential introns. Possibly, cross-correlations from two colored probes both inserted closer to the 3' end could be optimal designs.

Our method requires knowing the design of the experimental system (number and position of the loops), the speed of polymerase as input and calibrating the maximal fluorescence from

one gene. While the polymerase speed is an important parameter and erroneous assumption could influence the inference, we have shown that our inference is relatively insensitive to polymerase speeds (see S11 Fig). In the current experiments we do not have an independent calibration of the maximal fluorescence coming from one gene, which could introduce potential errors in our analysis. However the reproducibility of our results suggests that these potential errors are small.

We assumed an effective model that describes the transcription state of the whole gene and does not explicitly take into account the individual binding sites. As a result all the parameters we learn are effective and describe the overall change in the expression state of the gene and not the binding and unbinding of Bicoid to the individual binding sites. For concreteness we presented our model assuming a change in the promoter state and constitutive polymerase binding, but our current model does not discriminate between situations where the transcriptional kinetics are driven by polymerase binding and unbinding and promoter kinetics. The presented formalism can be extended to more complex scenarios that describe the kinetics of the individual binding sites and random polymerase arrival times. Since we already have little resolution power to discriminate between these effective models, we chose to interpret the results of only these effective models. The exact contribution of the individual transcription binding sites could be inferred from the activity of promoters with mutated binding sites. Similarly, other more complex models, such as a reversible three state model, or a model with many ON states, have not been ruled out by our current analysis but are possible within the current framework.

The time traces we had to analyze are very short and finite size effects are pronounced. Unlike in cell culture studies, where long time traces are available, we could not collect enough ON and OFF time statistics to characterize the promoter dynamics from the waiting time distributions. In this paper we show that simple statistics, the auto- and cross-correlation functions are powerful general tools that can be used in these kinds of challenging circumstances. To reach our final conclusion we had to combine different kinds of statistics, which is also a useful strategy when limited by data.

The approach we propose is a general method that can be used for any type of time trace analysis. However it becomes very useful when studying *in vivo* biological processes, where the biology naturally limits the available statistics. In our case the number of ON and OFF events is naturally limited by the short duration of the cell cycles. Our method explicitly calculates correlation functions for short traces, correcting for the finite size effects, and can be also used without making steady state assumptions about the dynamics (although this requires collecting sufficient statistics about two time points, which may be hard for short traces). With these corrections we see that while an effective two state model of the underlying dynamics of transcription regulation holds in the anterior and boundary regions of the embryo in all of the early cell cycles, the rates are different in the boundary and anterior regions, showing a strong dependence on position dependent factors such as Bicoid or maternal and zygotic Hunchback concentrations. More statistics will make it possible to build more explicit models of Bicoid dependent activation.

While our method is able to deconvolute the effects of the fluorescent probe and infer rates below the buffering limit of the probe (in our case $\tau_{buff} \sim 72$ s, see Fig 3F), in all cases, the rates that we can infer from time dependent traces are naturally limited by the timescales at which the polymerase leaves the promoter, which in our case is estimated to be $\tau_{block} \sim 6$ s. If the switching rates are faster than this scale, even a perfect, noiseless and infinitely accurate sampling of the dynamics will not be able to overcome this natural limit.

The inferred rates are reproducible between nuclei and embryos and the inter-embryo variability is similar to the inner-embryo variability (Fig 5A, 5C and 5D). The embryo-to-embryo variability can come from Bicoid variability, which is $\sim 10\%$ [27], so we do not expect the observed expression variability to be less, variability in growth rate and RNAP availability and external environmental factors. Additional sources of noise are experimental noise and most importantly problems with data calibration of what is the maximal level of fluorescence intensity.

We used the obtained results to estimate the precision of the transcriptional process from the *hunchback* promoter. We found that even in the anterior region, the variability in the mRNA produced in steady state by the different nuclei is large, with a relative error of about 50% (Fig 7A). This variability further increases to 150% of the mean mRNA produced at the boundary. These empirical estimates are completely explained for a two state promoter model by theoretical arguments, which treat the gene as an independent measuring device that samples the environment, correcting for the number of independent measurements during a cell cycle. In both cases, the precision at the level of the gene readout is not sufficient to form the precise Hunchback boundary up to half a nuclear width [39]. Even extending our argument to the total mRNA produced in the early cell cycles (Fig 7B) does not help. Having an irreversible promoter cycle could increase the theoretical precision, but only slightly in the parameter regime we have inferred and it would not change the quantitative conclusions about low precision backed by the empirical results. A Poisson-like promoter, while not compatible with the observed error rates, does have a significantly smaller error.

The construct we used here was limited to the 500 bp of the proximal hunchback promoter, which is known to recapitulate the hunchback endogenous expression observed in Fluorescent In Situ Hybridization (FISH) [20]. It is possible that the boundary phenotype is recovered by averaging of mRNAs and proteins produced by the real gene or the transgenes in other nuclei. In the latter case, this would point towards a robust "safety" averaging mechanism that relies on the population. Alternatively, we have to be aware that the sharp boundaries were only detected on fixed samples and that having access to the dynamics of the transcription process likely provides a more accurate view on the process. We calculated and estimated from the data the precision of the gene readout based on the variability of the transcription process between nuclei. We find that the transcriptional process at a given position is quite noisy. Previous estimates of precision were based on data from fixed samples and did not consider the probability of the gene to be ON, but assumed a binary representation where each nuclei is either active or inactive. By analyzing the full dynamic process we show that the gene is bursty and the transcriptional process itself is much more variable. Reducing the information contained in our traces to binary states, we find precise expression in the anterior, but still large variability at the boundary, similarly to previous results from Fluorescent In Situ Hybridization (FISH) aiming to detect all mRNAs [6].

Assuming that the precision in determining the position of the nuclei is encoded in the precision of the gene readout, a gene with the dynamics characterized in this paper needs to measure the signal ~ 200 times longer at the boundary to achieve the observed $\sim 10\%$ precision. A gene in the anterior would need to integrate only ~ 25 times longer. These results again suggest that the precision in determining the position of the nuclei is not only encoded in the time averaged gene readout, but probably relies either on spatial averaging mechanisms [27, 29, 40] or more detailed features of the temporal information encoded in the full trace [32].

In summary, the early developing fly embryo provides a natural system where we can investigate a functional setting the dynamics of transcription in a living organism. In our data analysis we are confronted with the same limitations that natural genes face: an estimate of the environmental conditions must be made in a very short time. Analysis of dynamical traces suggests that transcription is a bursty process with relatively large inter-nuclei variability, suggesting that simply the templated one to one time-averaged readout of the Bicoid gradient is unlikely. Comparing mutant experiments can shed light on exactly how the decision to form the sharp *hunchback* mRNA and protein boundary is made.

Materials and Methods

Experimental procedures

Constructs. For live monitoring of *hb* transcription activity in Drosophila embryos, we used the MS2-MCP system, which allows fluorescent labeling of RNAs as they are being transcribed [17, 38, 41]. To implement the reporter system in embryos, we generated flies transgenic for single insertions of a P-element carrying the *hb* proximal promoter upstream of an iRFP-MS2 cassette carrying 24 MS2 repeats [17, 42], from which Zelda binding motifs have been removed [20]. The flies also carry the P{mRFP-Nup107.K} [43] transgenic insertion on the 2nd chromosome and the Pw[+mC] = Hsp83-MCP-GFP transgenic insertion on the 3rd chromosome. These allow the expression of the Nucleoporin-mRFP (mRFP-Nup) for the labeling of the nuclear envelopes and the MCP-GFP required for labeling of nascent RNAs [41]. All stocks were maintained at 25°C.

Live imaging. Embryo collection, dechorionation and imaging have been done as described in [17]. Image stacks ($\sim 19Z \times 0.5\mu$ m, 2μ m pinhole) were collected continuously at 0.197μ m XY resolution, 8bits per pixels, 1200x1200 pixels per frame. A total of 4 movies capturing 4 embryos from nuclear cycle 10 to nuclear cycle 13 were taken. Each movie had different scanned fields along the embryo's width, which results in different time resolutions: 13.1 *s*, 10.2 *s*, 5.1 *s* and 4.3 *s*.

Image analysis. Nuclei segmentation, tracking and MS2-MCP loci analysis were performed as in [17] and recapitulated here. All steps were inspected visually and manually corrected when necessary. Nuclei segmentation and tracking were done by analyzing, frame by frame, the maximal Z-projection of the movies' mRFP-Nup channel. Each image was fit with a set of nuclei templates, disks of adjustable radius and brightness comparable with those of raw nuclei, from which the nuclei's positions were extracted. During the cycle's interphase, each nucleus was tracked over time with a simple minimal distance criterion. For MS2-MCP loci detection and fluorescent intensity quantification, the 3D GFP channel (MS2-MCP) were masked with the segmented nuclear images obtained in the previous step. This procedure also helps associating spots to nuclei. We then applied a threshold equal to ~ 2 times the background signal to the masked images and selected only the connected regions with an area larger than 10 pixels. The spot positions were set as the position of the centroids of the connected regions. The intensity of each spot was calculated by summing up the total pixel intensity in the vicinity of the centroids (a region of 1.5μ m x 1.5μ m x 1μ m) subtracted from the background intensity, which was extracted from the nearby region but excluding the spots. The flies were heterozygous for the MS2 reporter, so one spot was visible at a time. In the (rare) case of multiple spots detected per nucleus, the biggest spot was selected.

For each nucleus, we collected the nucleus' position and the spot intensity over time (here referred to as "traces"). The traces were then classified according to their respective embryos (out of 4 embryos), cell cycle (10 to 13) and position along the AP axis (either anterior or boundary). See SI Section A and S1 Fig for examples of traces.

Trace preprocessing. Before the autocorrelation function can be calculated the traces need to be preprocessed. To ensure that the data captures the dynamics of gene expression in its steady state, for each embryos and each cell cycle, we observed the spot intensity only in a specific time window. The beginning and the end of this window is determined as the moment the mean spot intensity over time of all traces (both at the anterior and the boundary) reaches and leaves an expression plateau (see example in S2 Fig).

The two state model. The detailed form of the autocorrelation function in Eq.2 depends on the underlying gene promoter switching model. For the two state—telegraph switching model (left panel in Fig 1A) the jumping times between the two states are both exponential and the dynamics is Markovian. The mean steady state probability for the promoter to be ON is $P_{\text{on}} = k_{\text{on}}/(k_{\text{on}} + k_{\text{off}})$, which combined with Eq.1 gives the form of the mean fluorescence $\langle F \rangle$. The probability that the gene is ON at time *n* given that it was on at time 0 is $A_n = P_{\text{on}} + e^{(\delta-1)n}(1 - P_{\text{on}})$, where $\delta = 1 - k_{\text{on}} - k_{\text{off}}$. The steady state connected correlation function depends only on the time difference (see SI Section B):

$$\tilde{C}_{\tau} = \langle F(t)F(t+\tau) \rangle - \langle F(t) \rangle^2 = \sum_{i,j} L_i L_j P_{\rm on}(1-P_{\rm on}) e^{(\delta-1)|\tau-j+i|}.$$
(6)

The cycle model. In the cycle model (also called the three state model in the text) (center panel in Fig 1A) the OFF period is divided into different sub-steps that correspond to K intermediate states with exponentially distributed jumping times from one to the next. The transition matrix T encodes the rates of this irreversible chain. The probability of the promoter to be in the ON state is:

$$P_{\rm on} = \frac{k_{\rm off}^{-1}}{k_{\rm off}^{-1} + \sum_{m=1}^{K} k_m^{-1}},\tag{7}$$

and that the steady state connected autocorrelation function is (see SI Section D):

$$\tilde{C}_{\tau} = \langle F(t)F(t+\tau) \rangle - \langle F(t) \rangle^{2}
= \sum_{i=1}^{r} \sum_{j=1}^{r} L_{i}L_{j}P_{\text{on}} \begin{bmatrix} (1 \quad 0 \quad \dots \quad 0)e^{(T-\text{Id})*|i-j-\tau|} \begin{pmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix} - P_{\text{on}} \end{bmatrix},$$
(8)

where τ is counted in polymerase steps, Id is the identity matrix of dimension K + 1, and the two unit vectors are of dimension K + 1. In the simple case of a two state model Eq.8 reduces to Eq.6.

The Γ waiting time model. An alternative description of a promoter cycle relies on a reduced description to an effective two state model where we use the fact that the transitions between the states are irreversible. The distribution of times spent in the effective OFF state τ , is no longer exponential, as it was in the two state model, but it has a peak at nonzero waiting times, which can be approximated by a Gamma distribution

$$\Gamma(\tau) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} x^{\alpha - 1} e^{-\beta x}, \qquad (9)$$

with mean α/β , where β is the scale parameter, α is the shape parameter and $\Gamma(\alpha)$ is the gamma function. The true distribution of waiting times in a cycle model approaches the Γ distribution if the rates k_i are all the same and $k_i <<1$. In this limit $\beta \approx k_i$, and α describes the number of intermediate OFF states. In the more general case it correctly captures the effective properties of the process. The mean probability of the promoter to be in the ON state in the Γ waiting

time model is given by

$$P_{\rm on} = (1 + \frac{\alpha k_{\rm off}}{\beta})^{-1}.$$
 (10)

The autocorrelation function cannot be computed directly analytically. The steady state Fourier transform of the steady state autocorrelation is (see SI Section E):

$$\mathcal{F}(\langle F(t)F(t+\tau)\rangle - \langle F(t)\rangle^2)(\xi) = \int_{-\infty}^{+\infty} d\tau e^{-2i\pi\tau} (\langle F(t)F(t+\tau)\rangle - \langle F(t)^2\rangle)$$

$$= \sum_{k,j} L_k L_j P_{\rm on} 2\Re \left[e^{-2i\pi(i-j)} \left[\left(k_{\rm off} + 2i\pi\xi - k_{\rm off} \left(1 + \frac{2i\pi\xi\beta}{\alpha k_{\rm off}} \right)^{-\alpha} \right)^{-1} - \frac{P_{\rm on}}{2i\pi\xi} \right] \right].$$
(11)

Finite cell cycle length correction to the connected autocorrelation function. Due to the short duration of the cell cycle, the theoretical connected correlation functions need to be corrected for finite size effects when comparing them to the empirically calculated correlation functions. When analyzing the data we calculate the autocorrelation function from *M* traces $\{v_{\alpha}\}_{1 \le \alpha \le M}$ of the same length *K*, $v_{\alpha} = \{v_{\alpha j}\}_{1 \le j \le K}$. We calculate the connected autocorrelation function for each trace and normalize it to 1 at the second time point to avoid spurious nucleus to nucleus variability:

$$c_{\alpha}(r) = \frac{\sum_{(i,j),|i-j|=r} \left\{ \left(v_{\alpha i} - \frac{1}{K} \sum_{l=1}^{K} v_{\alpha l} \right) \left(v_{\alpha j} - \frac{1}{K} \sum_{l=1}^{K} v_{\alpha l} \right) \right\}}{\frac{K - r}{K} \sum_{j=1}^{K} \left(v_{\alpha j} - \frac{1}{K} \sum_{l=1}^{K} v_{\alpha l} \right)^{2}},$$
(12)

and then average over all *M* traces to obtain the final connected autocorrelation function:

$$c(r) = \frac{1}{M} \sum_{\alpha=1}^{M} c_{\alpha}(r).$$
(13)

We define $\bar{\nu} = \langle \nu_i \rangle$ —the steady state true theoretical average of the random fluorescence intensity over random realizations of the process, and $\bar{\nu}^2 = \langle \nu_i^2 \rangle$ —the true theoretical second moment of the fluorescence signal. When $K \to \infty$ the average over time points is equal to the theoretical average, $1/K \sum_{i=1}^{K} \nu_{\alpha i} = \bar{\nu}$. Using time invariance in steady state the autocorrelation function becomes:

$$C_r = \frac{\langle v_i v_{i+r} \rangle - \bar{v}^2}{\bar{v}^2 - \bar{v}^2},\tag{14}$$

where $\langle \cdot \rangle$ is an average over random realizations of the process. Eq 14 corresponds to the limit we calculated in the theoretical models. To account for the finite size effects that arise due to short time traces, we need to correct for the fact that for short traces $\frac{1}{K} \sum_{i=1}^{K} v_{\alpha i} \neq \bar{v}$ and

 $\frac{1}{K}\sum_{i=1}^{K} v_{\alpha i}^2 \neq \bar{v^2}$, instead both the mean and the variance are functions of *K*. We note that for short

traces the definitions of autocorrelation and autocovariance differ:

$$\sum_{(i,j),|i-j|=r} \left\{ \left(v_{\alpha i} - \frac{1}{K} \sum_{l=1}^{K} v_{\alpha l} \right) \left(v_{\alpha j} - \frac{1}{K} \sum_{l=1}^{K} v_{\alpha l} \right) \right\} \neq \sum_{(i,j),|i-j|=r} \left(v_{\alpha i} v_{\alpha j} - \frac{1}{K^2} \sum_{l=1}^{K} v_{\alpha l} \sum_{m=1}^{K} v_{\alpha m} \right).$$
(15)

In practice for the analyzed dataset we found that the finite size effects for the variance can be neglected, however the mean over time points is a bad approximation to the ensemble mean. We present the finite size correction to the mean below. For completeness we include the finite size correction for the variance in SI Section I, although we do not use it in the analysis due to its numerical complexity and small effect.

If the variance of the normalized fluorescence intensity over random realizations of the process is well approximated by the average over the *K* time points, we can replace the denominator in Eq 12 by $\bar{\nu}^2 - \bar{\nu}^2$ and in steady state evaluate the mean connected autocorrelation function (see SI Section H for details):

$$c(r) = \frac{1}{\bar{v^2} - \bar{v}^2} \left[\tilde{C}_r + \frac{1}{K} \left(\frac{1}{K} - \frac{2}{(K-r)} \right) \left(K \tilde{C}_0 + \sum_{k=1}^{K-1} 2(K-k) \tilde{C}_k \right) + \frac{2}{K(K-r)} (r \tilde{C}_0 + \sum_{k=1}^{r-1} 2(r-k) \tilde{C}_k + \sum_{m=1}^{K-1} \tilde{C}_m[\min(m+r,K) - \max(r,m)]) \right] (16)$$

where $C_k = \langle v_i v_{i+k} \rangle$ is the infinite-size steady state non-connected correlation function of the process (given in Eq.6 for the two state model, Eq.8 for the cycle model and as the Fourier transform of Eq.12 for the Γ model) and the average is over random realizations of the process. The mean and variance of the signal, \bar{v} and $\bar{v^2}$, provide a normalization factor that is constant for all time differences *r*. We normalize the autocorrelation function setting the second term to 1 and these terms are not needed for the inference. If $v_i = X(i)$ then C_k is proportional to A(k).

Inference. The inference proceeds in three steps:

- Step 1. Signal calibration. The intensity of the measured signal depends on a constant tracedependent offset value I_0 , $I(t) = \sum_{i=1}^r I_0 a(i, t) L_i$. To calibrate this offset we take the maximum expression to be the mean of the maximum expression over all traces in a given region $I_{\text{max}} = \langle \max_i I(t) \rangle = I_0 \sum_{i=1}^r L_i$. We take the mean of the maxima of the intensities rather than the absolute maximum of all signals to avoid errors from overestimating the maximum. The calibrated fluorescence signal used in the analysis is then $F(t) = I(t)/I_0 = \sum_{i=1}^r a(i, t)L_i$. P_{on} is directly calculated using Eq 1. L_i is a known function.
- Step 2. Estimating parameter ratios. For the two state, three state cycle and Γ models, the ratios of the rates can be estimated directly from the steady state mean fluorescence values using Eqs 7 and 10. The Poisson model is uniquely defined by $P_{\rm on}$ and does not require further parameter inference beyond Step 1.
- Step 3. Estimating parameters. Using the estimate for the ratio of the rates, the ON and OFF rates are found by minimizing the mean squared error between the autocorrelation function calculated from the data (Eq 13), and the model (Eq 16 with the theoretical prediction for the appropriate model: Eq 6 for the two state, Eq 8 for the cycle model and the Fourier transform of Eq 12 for the Γ model).

Supporting Information

S1 Text. Supplementary materials and methods. (PDF)

S1 Fig. Examples of individual spot intensity over time. Consecutively shown are the traces in (A) Cycle 12, Anterior, (B) Cycle 12, Boundary (C) Cycle 13, Anterior, (D) Cycle 13, Boundary. The x axis is time in minute and y axis is the spot intensity in AU. (TIF)

S2 Fig. Data calibration. Shown are examples of 5 (out of 154) individual traces (blue) taken from embryo 1, cycle 13. Also shown is the mean spot intensity over time of all traces (red). The steady state window is chosen to be from the 6^{th} minute to the 11^{th} minute (dashed lines). (TIF)

S3 Fig. Fits of the autocorrelation function. The empirical autocorrelation function (blue dots) for both the anterior and boundary regions in all four embryos is fit using the autocorrelation function with the finite size corrections for the Poisson-like model (red lines), two-state model (green lines) and three-state cycle model (black lines). (TIF)

S4 Fig. Example of the connected autocorrelation function for the two state model calculated for different trace lengths T. The shaded areas denote the standard variation over 500 simulated traces. The switching rates $k_{on} = k_{off} = 0.01s^{-1}$. (TIF)

S5 Fig. The finite trace effect for the Ornstein-Uhlenbeck process. The connected autocorrelation function $C_r = \exp(-t/\tau)$ (red line) compared to the connected autocorrelation function calculated from short time traces as described in SI Section I H (blue line) and the corrected connected autocorrelation function (Eq. 54). $\lambda = 2s^{-1}$, $\gamma = 4s^{-1/2}$ and the short trace length is 5s where the Ornstein-Ulhenbeck process is $\partial_t x = -\lambda x + \gamma \xi$ and ξ is Gaussian white noise.

(TIF)

S6 Fig. Inference of the two-state model from the cross-correlation function between 3' signals and 5' signals. The gene cassette contains two identical arrays of MS2 binding sites on the 3' and 5' ends, separated by a gene of 3 kbp in length. The input parameters k_{on} , k_{off} are varied so as to maintain the same $P_{on} = 0.1$. (TIF)

S7 Fig. Comparison between the autocorrelation functions of the Poisson-like model and the two-state model. Shown are the autocorrelation functions (calculated from 1000 traces of 250 s in length) of the Poisson-like model (dashed black) and the two-state model (solid) with varying k_{on} and k_{off} . The model parameters are set to achieve the same effective transcription rate, $P_{on} = 0.1$, that we infer in the boundary region. For large $k_{on} + k_{off}$ values the shape of the autocorrelation function is dominated by the autocorrelation of the fluorescent probe and the Poisson-like and two state model autocorrelation functions look very similar. The inferred two state parameters are close to the green line. Since it is difficult to estimate the number of independent measurements, we cannot use standard statistical measures to compare these models with different numbers of parameters, whereas to determine the value of parameters within a given model we use a statistical measure (the mean square distance between the model prediction and data). For this reason we can differentiate between parameter values for the two state model that result in similar looking autocorrelation functions, but we cannot differentiate between two classes of models that result in similar differences in the autocorrelation functions.

(TIF)

S8 Fig. The fit of the three state cycle model to the data. The fit of the ratio of the two rates for leaving the two OFF states, k_1/k_2 , to the steady state traces from four embryos in the anterior and boundary region of cell cycle 12 and 13. Each point is data from one embryo. The error bars represent the standard deviation of the inferred value. The fit is for a randomized 60% of the data. The sum of the switching rates $k_{on} + k_1 + k_2$ is shown in Fig 5B of the main text.



S9 Fig. The relative error of gene expression. A. The mean probability of the gene to be ON at any time during the cell cycle as a function of the embryo length (binary approximation). B. Comparison of the relative error in the mRNA produced during the steady state of the interphase estimated empirically from data (abscissa) and from theoretical arguments in Eq. 62 using the inferred parameters from the autocorrelation function (ordinate), in the anterior (blue) and the boundary (red) regions, show very good agreement. C. The conclusions about precision do not depend on the embryo. The relative error of the total mRNA produced in cell cycle 13 as a function of position for windows equal to 10% of the embryo length. Each colored line represents one embryo. The same data plotted as an average over embryos with the variance as error bars is shown in Fig 7 of the main text. D. The conclusions about precision do not depend on the window size. The total mRNA produced in cell cycle 13 as a function of position for different window sizes. Except for very large scales (20%) and very small scales comparable to one nuclear width (2%, the relative error as a function of position is reproducible. E. The mean probability for the gene to be ON averaged over the cell cycle. F. The relative error of the discrete variable that describes the probability of the gene to be ON at any time during the cell cycle as function of position. The relative error is much lower in the anterior compared to the error in the total produced mRNA, but remains high at the boundary. In A, C, E and F each colored lines describe different embryos.

(TIF)

S10 Fig. The autocorrelation function for Poisson-like model and the two-state model for infinitely-long time traces. Autocorrelation functions of the Poisson-like model (dashed black) and the two-state models (solid) with $P_{on} = 0.1$ (similar to the inferred value in the boundary region) and varying $k_{on} + k_{off}$. In the inferred parameter regime (approximately green line), longer time traces do not help distinguish the two models based on the autocorrelation function. For large $k_{on} + k_{off}$ values the shape of the autocorrelation function is dominated by the autocorrelation of the fluorescent probe and the Poisson-like and two state model autocorrelation functions look very similar, even for long traces. (TIF)

S11 Fig. The dependence of the data fit on polymerase blocking time. Assuming different buffering times for the polymerase does not strongly affect the fit of the switching rates: a fit with $\tau_{block} = 4s$ (A) and $\tau_{block} = 8s$. $\tau_{block} = 6s$ is used in the main text in Fig 5D. (TIF)

Acknowledgments

We thank T. Mora for helpful discussions.

Author Contributions

Conceptualization: JD HT TF TL CPR CF MC ND AMW.
Data curation: JD HT TF TL CPR CF MC ND AMW.
Formal analysis: JD HT TF TL CPR CF MC ND AMW.
Funding acquisition: CF MC ND AMW.
Investigation: JD HT TF TL CPR AG CF MC ND AMW.
Methodology: JD HT TF TL CPR AG CF MC ND AMW.
Project administration: JD HT TF TL CPR AG CF MC ND AMW.
Resources: JD HT TF TL CPR AG CF MC ND AMW.
Software: JD HT TF TL CPR CF MC ND AMW.
Supervision: JD HT TF TL CPR CF MC ND AMW.
Validation: JD HT TF TL CPR AG CF MC ND AMW.
Visualization: JD HT TF TL CPR CF MC ND AMW.
Writing – original draft: JD HT TF TL CF MC ND AMW.
Writing – review & editing: JD HT CF MC ND AMW.

References

- 1. Jaeger J. The gap gene network. Cellular and Molecular Life Sciences. 2011; 68(2):243–74. doi: 10. 1007/s00018-010-0536-y PMID: 20927566
- 2. Gregor T, Garcia HG, Little SC. The embryo as a laboratory: quantifying transcription in Drosophila. Trends in Genetics. 2014; 30(8):364–75. doi: 10.1016/j.tig.2014.06.002 PMID: 25005921
- Driever W, Nuesslein-Volhard C. The bicoid Protein Determines Position in the Drosophila Embryo in a Concentration-Dependent Manner. Cell. 1988; 54:95–104. doi: <u>10.1016/0092-8674(88)90183-3</u> PMID: <u>3383245</u>
- Tikhonov M, Little SC, Gregor T. Only accessible information is useful: insights from patterning Subject Category. Royal Society Open Science. 2015; 2(150486). doi: 10.1098/rsos.150486 PMID: 26716005
- 5. Porcher A, Dostatni N. The bicoid morphogen system. Current Biology. 2010; 20(5):R249–54. doi: 10. 1016/j.cub.2010.01.026 PMID: 20219179
- Little SC, Tikhonov M, Gregor T. Precise developmental gene expression arises from globally stochastic transcriptional activity. Cell. 2013; 154(4):789–800. doi: 10.1016/j.cell.2013.07.025 PMID: 23953111
- 7. Elowitz MB, Levine AJ, Siggia ED, Swain PS. Stochastic gene expression in a single cell. Science. 2002; 297(5584):1183–6. doi: 10.1126/science.1070919 PMID: 12183631
- Ozbudak EM, Thattai M, Kurtser I, Grossman AD, van Oudenaarden A. Regulation of noise in the expression of a single gene. Nature Genetics. 2002; 31(1):69–73. doi: 10.1038/ng869 PMID: 11967532
- Raser J, O'Shea E. Control of stochasticity in eukaryotic gene expression. Science. 2004; 304:1811. doi: 10.1126/science.1098641 PMID: 15166317
- Crauk O, Dostatni N. Bicoid determines sharp and precise target gene expression in the Drosophila embryo. Current Biology. 2005; 15(21):1888–98. doi: 10.1016/j.cub.2005.09.046 PMID: 16271865
- Suter DM, Molina N, Gatfield D, Schneider K, Schibler U, Naef F. Mammalian genes are transcribed with widely different bursting kinetics. Science. 2011; 332(6028):472–474. doi: 10.1126/science. 1198817 PMID: 21415320
- Zoller B, Nicolas D, Molina N, Naef F. Structure of silent transcription intervals and noise characteristics of mammalian genes. Molecular Systems Biology. 2015; 11(7):823. doi: 10.15252/msb.20156257 PMID: 26215071

- Taniguchi Y, Choi PJ, Li GW, Chen H, Babu M, Hearn J, et al. Quantifying E. coli proteome and transcriptome with single-molecule sensitivity in single cells. Science. 2010; 329(5991):533–8. doi: 10. 1126/science.1188308 PMID: 20671182
- Kandhavelu M, Lloyd-Price J, Gupta A, Muthukrishnan AB, Yli-Harja O, Ribeiro AS. Regulation of mean and noise of the in vivo kinetics of transcription under the control of the lac/ara-1 promoter. FEBS Letters. 2012; 586(21):3870–5. doi: 10.1016/j.febslet.2012.09.014 PMID: 23017207
- Muthukrishnan AB, Kandhavelu M, Lloyd-Price J, Kudasov F, Chowdhury S, Yli-Harja O, et al. Dynamics of transcription driven by the tetA promoter, one event at a time, in live Escherichia coli cells. Nucleic Acids Research. 2012; 40(17):8472–83. doi: 10.1093/nar/gks583 PMID: 22730294
- Chong S, Chen C, Ge H, Xie XS. Mechanism of transcriptional bursting in bacteria. Cell. 2014; 158 (2):314–26. doi: 10.1016/j.cell.2014.05.038 PMID: 25036631
- Lucas T, Ferraro T, Roelens B, Chanes JDLH, Walczak A, Coppey M, et al. Live imaging of Bicoiddependent transcription in Drosophila embryos. Current Biology. 2013; 23(21):2135–2139. doi: 10. 1016/j.cub.2013.08.053 PMID: 24139736
- Garcia HG, Tikhonov M, Lin A, Gregor T. Quantitative imaging of transcription in living Drosophila embryos links polymerase activity to patterning. Current Biology. 2013; 23(21):2140–5. doi: 10.1016/j. cub.2013.08.054 PMID: 24139738
- Ferraro T, Lucas T, Clémot M, De Las Heras Chanes J, Desponds J, Coppey M, et al. New methods to image transcription in living fly embryos: the insights so far, and the prospects. Wiley interdisciplinary reviews: Developmental Biology. 2016;. doi: 10.1002/wdev.221 PMID: 26894441
- 20. Lucas T, et al. in preparation. 2016;.
- Bothma JP, Garcia HG, Esposito E, Schlissel G, Gregor T, Levine M. Dynamic regulation of eve stripe 2 expression reveals transcriptional bursts in living Drosophila embryos. Proceedings of the National Academy of Sciences of the United States of America. 2014; 111(29):10598–603. doi: <u>10.1073/pnas.</u> 1410022111 PMID: 24994903
- 22. Estrada J, Wong F, DePace A, Gunawardena J. Information Integration and Energy Expenditure in Gene Regulation. Cell. 2016; 166(1):234–44. doi: 10.1016/j.cell.2016.06.012 PMID: 27368104
- Coulon A, Ferguson ML, de Turris V, Palangat M, Chow CC, Larson DR. Kinetic competition during the transcription cycle results in stochastic RNA processing. eLife. 2014; 3:1–22. doi: <u>10.7554/eLife.03939</u> PMID: 25271374
- 24. Coulon A, Larson DR. Fluctuating Analysis: Dissecting Transcriptional Kinetics with Signal Theory. Methods in Enzymology. 2016; 03:1–33. doi: 10.1016/bs.mie.2016.03.017 PMID: 27241754
- Gregor T, Wieschaus EF, McGregor AP, Bialek W, Tank DW. Stability and nuclear dynamics of the bicoid morphogen gradient. Cell. 2007; 130(1):141–52. doi: 10.1016/j.cell.2007.05.026 PMID: 17632061
- Cheung D, Miles C, Kreitman M, Ma J. Adaptation of the length scale and amplitude of the Bicoid gradient profile to achieve robust patterning in abnormally large Drosophila melanogaster embryos. Development (Cambridge, England). 2014 Jan; 141(1):124–35. doi: 10.1242/dev.098640 PMID: 24284208
- 27. Gregor T, Tank DW, Wieschaus EF, Bialek W. Probing the limits to positional information. Cell. 2007; 130(1):153–64. doi: 10.1016/j.cell.2007.05.025 PMID: 17632062
- Berg HC, Purcell EM. Physics of chemoreception. Biophysical Journal. 1977; 20(2):193–219. doi: 10. 1016/S0006-3495(77)85544-6 PMID: 911982
- Erdmann T, Howard M, Ten Wolde PR. Role of spatial averaging in the precision of gene expression patterns. Physical Review Letters. 2009; 103(25):2–5. doi: 10.1103/PhysRevLett.103.258101 PMID: 20366291
- Porcher A, Dostatni N. The Bicoid morphogen system. Current Biology. 2010; 20(5):249–254. doi: 10. 1016/j.cub.2010.01.026
- Perry MW, Bothma JP, Luu RD, Levine M. Precision of hunchback expression in the Drosophila embryo. Current Biology. 2012; 22(23):2247–2252. doi: 10.1016/j.cub.2012.09.051 PMID: 23122844
- He F, Ren J, Wang W, Ma J. A multiscale investigation of bicoid-dependent transcriptional events in Drosophila embryos. PloS one. 2011 Jan; 6(4):e19122. doi: 10.1371/journal.pone.0019122 PMID: 21544208
- Fukaya T, Lim B, Levine M. Enhancer Control of Transcriptional Bursting. Cell. 2016;p. 1–11. doi: 10. 1016/j.cell.2016.05.025 PMID: 27293191
- Wang Y, Tu KC, Ong NP, Bassler BL, Wingreen NS. Protein-level fluctuation correlation at the microcolony level and its application to the Vibrio harveyi quorum-sensing circuit. Biophysical Journal. 2011; 100 (12):3045–53. doi: 10.1016/j.bpj.2011.05.006 PMID: 21689539

- Lin Y, Sohn CH, Dalal CK, Cai L, Elowitz MB. Combinatorial gene regulation by modulation of relative pulse timing. Nature. 2015; 527(7576):54–8. doi: 10.1038/nature15710 PMID: 26466562
- Dunlop MJ, Cox RS, Levine JH, Murray RM, Elowitz MB. Regulatory activity revealed by dynamic correlations in gene expression noise. Nature Genetics. 2008; 40(12):1493–8. doi: <u>10.1038/ng.281</u> PMID: <u>19029898</u>
- Munsky B, Neuert G, van Oudenaarden A. Using gene expression noise to understand gene regulation. Science. 2012; 336(6078):183–7. doi: 10.1126/science.1216379 PMID: 22499939
- Ferraro T, Esposito E, Mancini L, Ng S, Lucas T, Coppey M, et al. Transcriptional Memory in the Drosophila Embryo. Current Biology. 2016; 26(2):212–8. doi: 10.1016/j.cub.2015.11.058 PMID: 26748851
- 39. Gregor T, Tkacik G, et al. oral communication. 2016;.
- 40. Okabe-Oho Y, Murakami H, Oho S, Sasai M. Stable, precise, and reproducible patterning of bicoid and hunchback molecules in the early Drosophila embryo. PLoS Computational Biology. 2009; 5(8): e1000486. doi: 10.1371/journal.pcbi.1000486 PMID: 19714200
- 41. Bertrand E, Chartrand P, Schaefer M, Shenoy SM, Singer RH, Long RM. Localization of ASH1 mRNA particles in living yeast. Molecular Cell. 1998; 2(4):437–445. doi: <u>10.1016/S1097-2765(00)80143-4</u> PMID: <u>9809065</u>
- Janicki SM, Tsukamoto T, Salghetti SE, Tansey WP, Sachidanandam R, Prasanth KV, et al. From silencing to gene expression: Real-time analysis in single cells. Cell. 2004; 116(5):683–698. doi: 10. 1016/S0092-8674(04)00171-0 PMID: 15006351
- Katsani KR, Karess RE, Dostatni N, Doye V. In Vivo Dynamics of Drosophila Nuclear Envelope Components. Molecular Biology of the Cell. 2008; 19:3652–3666. doi: <u>10.1091/mbc.E07-11-1162</u> PMID: 18562695