

RESEARCH ARTICLE

Approximated Gene Expression Trajectories (AGETs) for Gene Regulatory Network Inference on Cell Tracks

Kay Spiess^{1,2}, Shannon E. Taylor³, Timothy Fulton¹, Kane Toh¹, Dillan Saunders¹, Seongwon Hwang¹, Yuxuan Wang¹, Brooks Paige^{2,4*}, Ben Steventon^{1*} and Berta Verd^{1,3*}

ABSTRACT

The study of pattern formation has greatly benefited from our ability to reverse-engineer gene regulatory network (GRN) structure from spatio-temporal quantitative gene expression data. Traditional approaches omit tissue morphogenesis, and focus on systems where the timescales of pattern formation and morphogenesis can be separated. In such systems, pattern forms as an emergent property of the underlying GRN and mechanistic insight can be obtained from the GRNs alone. However, this is not the case in most animal patterning systems, where patterning and morphogenesis are co-occurring and tightly linked. To address the mechanisms driving pattern formation in such systems we need to adapt our GRN inference methodologies to explicitly accommodate cell movements and tissue shape changes. In this work we present a novel framework to reverse-engineer GRNs underlying pattern formation in tissues undergoing morphogenetic changes and cell rearrangements. By integrating quantitative data from live and fixed embryos, we approximate gene expression trajectories (AGETs) in single cells and use a subset to reverse-engineer candidate GRNs using a Markov Chain Monte Carlo approach. GRN fit is assessed by simulating on cell tracks (live-modelling) and comparing the output to quantitative data-sets. This framework generates candidate GRNs that recapitulate pattern formation at the level of the tissue and the single cell. To our knowledge, this inference methodology is the first to integrate cell movements and gene expression data, making it possible to reverse-engineer GRNs patterning tissues undergoing morphogenetic changes.

KEYWORDS: gene regulatory networks, network inference, pattern formation

INTRODUCTION

Embryonic pattern formation underlies much of the diversity of form observed in nature. As such, one of the main goals in developmental biology is to understand how spatio-temporal molecular patterns emerge in developing embryos, how they are maintained and how they can change over the course of evolution. Over the past three decades, the field has focused on the function and dynamics of the gene regulatory networks (GRNs) underlying these processes. GRNs can be formulated mathematically as non-linear systems of coupled differential equations whose parameters can be inferred from quantitative gene expression data: a methodology known as reverse-engineering (1; 2; 3; 4; 5; 6; 7; 8). Reverse-engineering has been successfully applied to a myriad of systems, from the *Drosophila* blastoderm to the vertebrate neural tube (9; 10; 11; 12), uncovering the mechanisms by which GRNs read out morphogen gradients (13; 14; 15; 16; 17), scale patterns (18), control the timing of differentiation (19; 20; 21), synchronise cellular fates (22) and evolve pattern formation (23).

Much of what we know about pattern formation has been learnt from reverse-engineering GRN structure from spatio-temporal quantitative data in systems where the timescales of pattern formation and morphogenesis are different and can therefore be separated. In such systems, spatio-temporal gene expression profiles are typically obtained by measuring gene expression levels across the tissue of interest in fixed stained samples, and interpolating between measurements at different time points (8). The underlying and seldom stated assumption, is that the gene expression dynamics are much faster than the cell movements in the developing tissue, and that therefore cell movements can be ignored over the timescales at which the pattern forms. This is true in many systems and processes, such as segmental patterning in early *Drosophila* embryogenesis. In systems where this is indeed the case, pattern formation can be considered an emergent property of GRN dynamics alone (24) and much insight can be drawn from analysing reverse-engineered GRNs (10; 13).

In systems where tissue patterning and tissue morphogenesis are coupled and occurring simultaneously, GRNs alone cannot account for the resulting patterns. This has been recently highlighted by work in organoids, where shape, size and cell type distribution are difficult to control as a result of altered patterning due to abnormal morphogeneses in unconstrained tissue geometries (25). Therefore, in order to be able to understand developmental pattern formation in a broader range of systems, we have to address how morphogenesis and GRNs together control fate specification and embryonic organisation. Importantly, to be able to do this, we need novel reverse-engineering methodologies that will explicitly accommodate cell movements and tissue shape changes.

¹Department of Genetics, University of Cambridge, Cambridge, UK;

²The Alan Turing Institute, London, UK;

³Department of Biology, University of Oxford, UK;

⁴Centre for Artificial Intelligence, University College London, UK

Authors for correspondence: berta.verdfernandez@biology.ox.ac.uk, bjs57@cam.ac.uk, b.paige@ucl.ac.uk

In this work we present a methodology to reverse-engineer GRNs underlying pattern formation in tissues that are undergoing morphogenetic changes such as cell rearrangements. As a case study we focus on T-box gene patterning in the developing zebrafish presomitic mesoderm (PSM) (Fig.1A). T-box genes coordinate fate specification along the PSM as cells move out of the tailbud and make their way towards the somites (26). Cell movements in the PSM can be live-imaged and followed in 3D (27). By the time they reach a somite, cells in the PSM will have undergone a stereotypical progression of T-box gene expression: *Tbx16* and *Tbx16* in the tailbud, followed by *Tbx16* in the posterior PSM and *Tbx6* in the anterior PSM (Fig.1A&D). The *Tbx16/Tbx6* boundary roughly marks the cells' transition out of the tailbud and in zebrafish it is thought to correlate with marked changes in cell behaviours where extensive cell mixing in the tailbud gives way to reduced, almost nonexistent mixing and neighbourhood cohesion in the PSM (28). Therefore, while all cells will eventually have undergone the same gene expression progression, their expression dynamics will differ as cells spend variable amounts of time in the tailbud (26). Despite this, a tissue-level pattern forms which scales with PSM length during the course of posterior development and somitogenesis (26). T-box pattern formation in the developing zebrafish PSM is therefore a good example of a developmental process where the molecular pattern across the tissue is an emergent property of the GRN, the cell movements and tissue shape changes involved in the tissue's morphogenesis.

The reverse-engineering methodology presented in this paper accommodates cell movements and tissue shape changes, representing tissue morphogenesis explicitly when reverse-engineering GRNs. To do this, our methodology integrates two different kinds of quantitative data: cell tracking data obtained from live-imaging the developing tissue and three-dimensional quantitative gene expression of the genes and signalling pathways of interest over developmental time. We project the 3D gene expression data onto the cell tracks to approximate gene expression trajectories (AGETs) in single cells. Using a subset of AGETs from ten cells pseudo-randomly spaced within the tissue we were able to reverse-engineer candidate GRNs applying a Markov Chain Monte Carlo (MCMC) approach. The fit of the resulting candidate GRNs is assessed by simulating them in each cell in the tracks using initial and boundary conditions extracted directly from the gene expression data, a methodology that we refer to as "live-modelling". The resulting well-fitting GRNs were grouped into 22 clusters, generating candidate GRNs that can be further investigated and challenged using experimental work (26).

To our knowledge, this inference methodology is the first to integrate cell movements and gene expression data, making it possible to reverse-engineer GRNs patterning tissues as they undergo morphogenesis. We hope that this toolbox will contribute to broaden the types of patterning systems that are studied quantitatively and mechanistically, increasing our understanding of pattern formation in development and evolution.

RESULTS AND DISCUSSION

Approximating gene expression dynamics on single cell tracks: AGETs

The ideal data to reverse-engineer gene regulatory networks would be temporally accurate quantifications of gene expression dynamics at the single cell level as the tissue develops. Unfortunately, current state of the art in live gene expression reporter technology,

while very advanced, cannot follow three genes and two signalling pathways simultaneously in space and time, while also ensuring that the dynamics of all reporters faithfully recapitulate the expression dynamics of the genes of interest. For this reason, it has been necessary to develop an alternative approach to effectively construct in-silico reporters which is based on approximating gene expression trajectories in the cells of the developing PSM, which we will from now on refer to as AGETs (approximated gene expression trajectories).

In brief, AGETs are obtained by projecting 3D spatial quantifications of gene expression in PSM cells obtained using HCRs and antibody stains, onto the cells present at each time frame of a time lapse of the developing PSM at approximately the same stages. The projected expression levels are used to assign gene and signalling expression levels in every cell in the time lapse. The result is an approximated gene expression trajectory for every cell in the time lapse, which can now be used to reverse-engineer gene regulatory networks which, when simulated on the tracks recapitulate T-box pattern formation on the developing PSM.

Data requirements and preparation

Two kinds of data are required to produce AGETs: cell tracks obtained from live-imaging the developing tissue of interest and quantitative spatial gene expression data at each developmental stage covered by the tracks.

In this case study, cell tracks were obtained by live-imaging a fluorescently labelled developing zebrafish tailbud between the 22nd and 25th somite stages using a two-photon microscope (see (27) and Materials and Methods). The raw data obtained consists of a series of point clouds representing the position of single cells in 3D space at 61 consecutive frames, which were taken at two minute intervals. The raw data were processed using a tracking algorithm in the image analysis software Imaris to obtain the position of single cells over time, and selected tracks were validated manually. The resulting data are a collection of cell tracks that describe the how individual cells move as the zebrafish tailbud and PSM develop. A cell track provides spatial information over time but is devoid of any information regarding gene expression levels in each cell.

Gene expression levels were approximated from fixed tailbud samples stained for the T-box gene products using HCR (29) and antibody stains for the signals Wnt and FGF (see Materials and Methods). If gene expression patterns don't scale with the development of the tissue, stage-specific stains should be used separately. Otherwise, if the pattern of interest scales with tissue growth over developmental time - as is the case in the developing zebrafish PSM - images at different, but close, stages can be quantified and pooled together. T-box genes - *Tbx16*, *Tbx16* and *Tbx6* - were simultaneously stained for on zebrafish tailbuds that had been fixed between the 23rd and 25th somite stages (SS) (Fig.1A). Of a total of 13 images, ten were processed and used for fitting (2x 23SS, 3x 24SS and 5x 25SS). Three separate antibody stained samples were used to quantify signals Wnt and FGF. In addition to the gene expression, tailbuds were stained with DAPI to be able to locate the cells by the position of their nuclei. Only one side of the zebrafish PSM was used.

A processing pipeline was developed to quantify the imaging data using the image analysis software Imaris (Fig.1). The first step in the pipeline consists of masking the PSM from the surrounding tissues, including the spinal cord and the notochord. This

was achieved by drawing a surface around the PSM using morphological and gene expression landmarks as a guide to identify different tissue boundaries (Fig.1B). Next, in order to consider only gene expression levels inside of the isolated PSM, all gene expression outside of the defined surface was set to zero (Fig.1C). Background noise in the data was reduced by setting lower-bound thresholds for every gene. These thresholds were chosen such that *Tbx16* and *Tbx16* would appear restricted to the posterior end of the PSM (Fig.1Di, Dii, Ei and Eii) with their expression in the anterior PSM reduced to zero. Similarly, thresholds were set for *Tbx6* expression to eliminate any background expression in the posterior PSM (Fig.1Diii and Eiii). Each gene is then normalized; normalization had to be robust enough to noisy gene expression levels. A Savitzky-Golay filter was applied to each gene to smoothen the signal (Fig.1D) and the smoothened maximum for each gene was set to one. Finally, spots were created in each detected nucleus from which a point cloud consisting of the 3D spatial coordinates and associated *Tbx16*, *Tbx16* and *Tbx6* levels were extracted (Fig.1E). The same pipeline was used to obtain the levels of signals Wnt and FGF in single cells.

AGET construction

AGETs are constructed to approximate the gene expression dynamics of single cells as they move and undergo complex re-arrangements during tissue morphogenesis. This requires live-imaging data, which provides information of the cell's spatial trajectories over time, to be combined with quantitative single cell gene expression data. To achieve this, we project the pre-processed HCR data (Fig.1E) onto the tracks to obtain an approximated read-out of the gene expression and signalling levels that each cell experiences as it moves.

The first step to project the extracted quantitative gene expression data onto the cell tracks is to align the point clouds representing the positions of the cells in 3D space processed from the HCRs (Fig.1E) with the point clouds for each of the 61 time frames in the time lapse (Fig.2A). We use point-to-plane ICP (iterative closest point) to perform this alignment (30), which in brief, is an iterative algorithm that seeks to map two point clouds onto each other by recursively minimising the distance between them (see Materials and Methods, and Fig.2). Once the point clouds have been aligned, equivalent regions of different PSMs will overlap in space (Fig.2A) making it possible to use the quantitative gene expression from cells in the processed HCRs to assign gene expression values to the cells (represented by points) in the time lapse at each time frame (Fig.2Bi and Bii and Algorithm1).

To approximate the gene expression and signalling values in a cell from the time lapse, we first find its five closest neighbouring cells from the processed HCR data. Since all PSMs have been aligned as point clouds, we now have a point cloud representing cells from both the PSM in the time lapse and those from the HCRs. The median gene expression and signalling values are calculated from the expression and signalling values of the five nearest neighbouring cells and assigned to the cell from the time lapse (Fig.2B and Algorithm 1 for a more detailed description of the process). Fig.2Bi shows the result of mapping T-box gene expression data from ten pre-processed HCR images onto the first frame of the tracking data and Fig.2Bii shows a quantification of the gene expression levels for all cells along the posterior to anterior axis. We repeat this procedure for each of the 61 frames in the time lapse resulting in an approximated gene expression trajectory (AGET) for every cell in the timelapse (Fig.2C and Supplementary

Movies 1 and 2). In addition, AGET construction was found to be robust to the specific number of neighbours used as well as to the method used to assign expression values at each time point (Supplementary Figures1 and 2)

Using AGETs to reverse-engineer gene regulatory networks that recapitulate pattern formation on a developing tissue

GRN models are often formulated as systems of coupled differential equations where state variables describe the concentrations of the gene products of interest and parameters represent the interactions between genes, as well as other factors such as production and degradation rates. In the case of the T-box genes, there are three state variables representing *Tbx16*, *Tbx16* and *Tbx6* levels and a total of 24 parameters to be fit (see Materials and Methods). Dynamic data are required to constrain and fit such models, and in this case these will be provided by the AGETs calculated previously. AGETs will be used as the target expression dynamics for the fitting procedure, where previously directly measured gene expression dynamics would have been used. As with other fitting procedures, an optimal parameter set will be one that minimises the difference between the target and the simulated data. We chose to adapt a Markov Chain Monte Carlo (MCMC) algorithm to use as our parameter sampling method since MCMC has been extensively used and repeatedly validated for GRN inference (31). In addition, MCMC has the advantage of providing a population of candidate networks by approximating the entire posterior distribution for each GRN parameter.

Using all 1903 available AGETs to fit our models would be ideal, as together they represent the tissue scale patterning dynamics that we seek to recapitulate. However, this is currently computationally expensive and ultimately unfeasible. Instead, we find that good fits are obtained when fitting to an ensemble of as few as ten AGETs provided that these span the length of the PSM. The ten AGETs were selected semi-randomly, where a randomly chosen set of ten AGETs would be visually inspected to ensure that they included cells distributed across the antero-posterior length of the PSM, and would otherwise be discarded. In addition, we only selected AGETs of maximal duration, namely those that corresponded to cells that had been consecutively tracked for the entire duration of the time lapse (61 frames). The ten AGETs used for reverse engineering and their approximate position in an idealised PSM are shown in Fig.2C.ii. For this case study, we found that reverse-engineering using ten AGETs generated well-fitting candidate networks while avoiding over-fitting and optimising the computational time required, however we also found that increasing the number of AGETs for reverse-engineering increased the proportion of networks producing high quality fits (see Materials and Methods, and Supplementary Figures 3 and 4). We expect that the specific number of AGETs required to obtain good fits will be problem-specific.

MCMC inference yields a collection of parameter sets or combinations (samples) that together approximate the posterior distribution of the GRN's parameters: for every parameter, we obtain a probability distribution across its values, which provides information about the values that are most likely to produce good fits. We first chose to explore the network corresponding to the parameter set with the overall highest posterior probability score: the maximum a posteriori - or MAP - sample (Fig.3A). We simulated each of the ten AGETs used during the fitting procedure

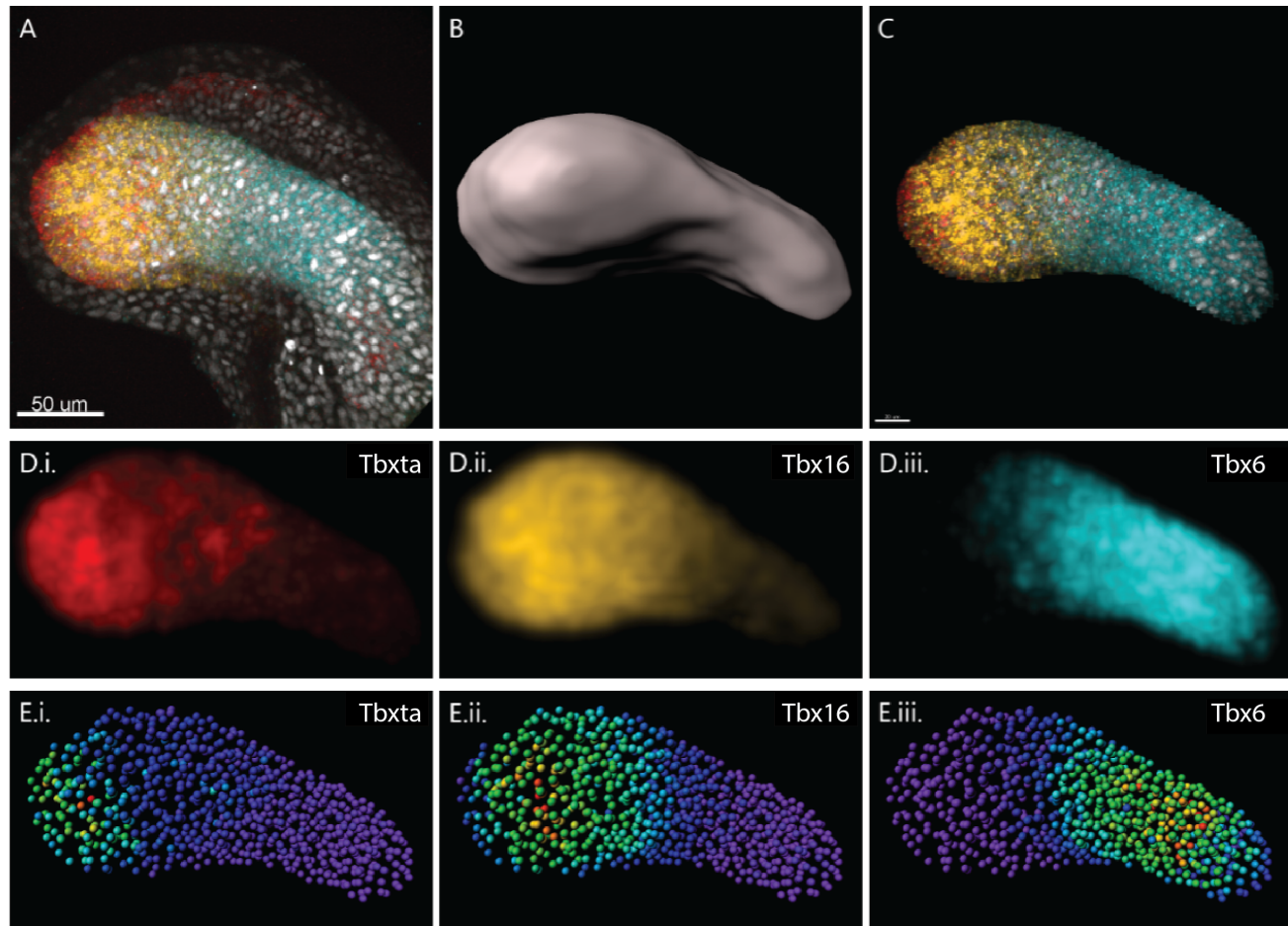


Fig. 1. Gene expression data preparation pipeline (A) Typical HCR image of a 22 somite stage zebrafish embryo tailbud stained for Tbxta (red), Tbx16 (yellow), Tbx6 (blue) and DAPI (gray). Anterior to the right, posterior to the left, dorsal up and ventral down from here on. (B) Surface masking the PSM based on T-box gene expression and morphological landmarks. (C) Gene expression and nuclear marker in the masked PSM (as before Tbxta in red, Tbx16 in yellow, Tbx6 in blue and DAPI in gray). (D) Normalising gene expression levels: Tbxta and Tbx16 levels in the anterior PSM are normalised to zero while posterior PSM levels of Tbx6 are normalised to zero, to eliminate background expression. A Gaussian filter has been then applied to each T-box gene to smoothen gene expression across the PSM. (E) Nuclei are segmented using the DAPI channel creating spots in 3D space. Spots are coloured according to the median intensity of each gene (i) Tbxta, ii) Tbx16 and iii) Tbx6), where purple denotes zero expression and red 1, which is the highest expression. The spatial coordinates of the spots together with the median intensities were exported and used to generate the AGETs.

and then simulated all 1903 available AGETs, and visualised the simulation on the tracks (Supplementary Movies 3 and 4). We validate the quality of the inferred network by both comparing single AGETs with their simulated counterparts (Fig.3B), and by comparing the whole tissue-level gene expression profiles over time (Fig.3C). We are especially interested in how well the simulations recapitulate whole tissue patterning dynamics, as these result from simulating AGETs that had not been used for model training. We discard parameter sets that simulate clear pattern aberrations, and consider a good fit to be when the position of gene expression domain intersections does not differ by more than the inter-embryonic biological boundary range (<10% A-P position) in the simulated versus the approximated patterns (Fig.3C). While quantitative measures of the goodness of fit can be easily defined, such as comparing the log-likelihood between parameter sets or calculating least-squares measures, these don't necessarily reflect whether aspects of the pattern that are of notable biological importance are being captured, and were therefore not favoured in this part of the analysis.

Fig.3B.i compares four of the ten AGETs (relative positions shown in Fig.3B.ii) (solid lines) used for model fitting with the resulting simulations (dotted lines). The simulated expression recapitulates well the target expression for the AGETs. The model was formulated as a deterministic system without added stochasticity which explains the smoothness of the simulated curves, which nonetheless can be seen to recapitulate AGET gene expression levels and trends. Fig.3C shows simulated T-box expression for each cell along the normalized posterior to anterior axis of the PSM (dots). The simulated data have been fit at each separate time point by curves which are then normalised (dotted curves) and compared to the curves previously fit in the same way to all AGETs (shown as solid curves). A comparison between AGETs and simulations is shown at three different time points in Fig.3C (simulation outputs at 33%, 66% and 100% total time respectively). Importantly, the overall position of the domains is recapitulated and the position of domain intersections is within the preset biological range of 10% A-P position. The full simulations can be found in Supplementary Movies 3 and 4.

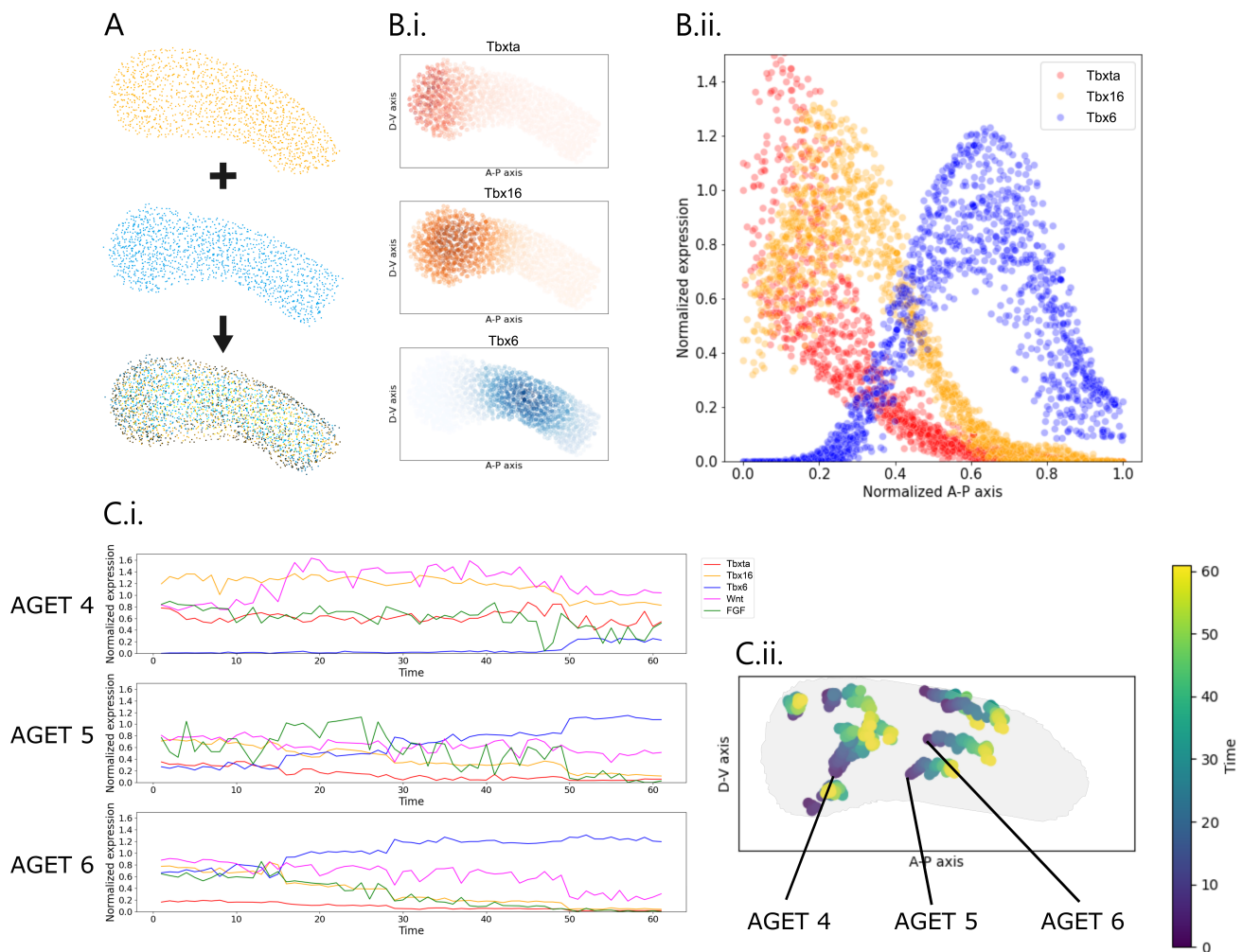


Fig. 2. Calculating AGETs. (A) In orange is the processed HCR image showing the positions of the cells in the PSM (source point cloud) and in blue are the positions of the cells taken from the first frame of the tracking data (target point cloud). Using ICP, all the source point clouds obtained from the HCR images are aligned with the target point cloud obtained from the first frame (61 in total) of the tracking data. This is illustrated by the overlapping orange, blue and black point clouds in the resulting point cloud (bottom). (B) i. T-box gene expression from ten pre-processed HCR images has been used to assign Tbx gene expression values to each cell in the first frame of the tracking data. Tbxta in red, Tbx16 in yellow and Tbx6 in blue. ii. Maximum projection of the data (first time point of the AGETs) in i. quantified along the posterior to anterior axis. (C) i. Three AGETs representing approximated T-box gene and signaling dynamics in three single cells at different position along the developing PSM (shown in C.ii). y-axis represents relative gene expression levels and x-axis reflects the time frame in the time lapse (from 1 to 61). Tbxta in red, Tbx16 in yellow and Tbx6 in blue, Wnt in pink, FGF in green. (C) ii. Ten cell tracks spanning the length of the PSM, whose AGETs were subsequently used for the GRN inference process. The ten cells have been chosen semi-randomly to cover the A-P axis. The outline illustrates the shape of the PSM. The color gradedness indicates time in timeframes. AGETs associated with cells 4, 5 and 6 are shown in panel C.i.

Notably, there is a discrepancy between the AGETs and the simulated anterior Tbx6 expression. The formulated GRN is unrealistic in this region, where additional factors secreted from the somites are known to be down-regulating this transcription factor (32). For this reason, it is reassuring and expected that the model doesn't recapitulate the pattern well in the anterior PSM border. In addition, the model predicts that over time, a small percentage of posterior cells will express low levels of Tbx6. Although unexpected, there is evidence suggesting that this is indeed the case (26). Such low and sparse posterior expression of Tbx6 would have been lost during the smoothing step in our data preparation pipeline, which is unable of capturing patterns of such fine resolution as it stands. It is encouraging that candidate GRNs consistently recapitulate this unexpected feature of the biology and

might suggest that the three genes considered are indeed causally responsible for most of this molecular patterning system.

Parameter determinability and model clustering

MCMC is a parameter sampling algorithm, and as such it will return an approximated posterior distribution for the GRN parameters instead of a single estimate. This provides a range of candidate networks that can be subsequently analysed and challenged in combination with experimental approaches. Such parameter distributions also provide valuable information regarding which model parameters — and therefore genetic interactions — are tightly constrained by the data, and which aren't and therefore appearing to take on a broad range of values across the inferred networks. Such information can lead to interesting hypotheses regarding which

561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616

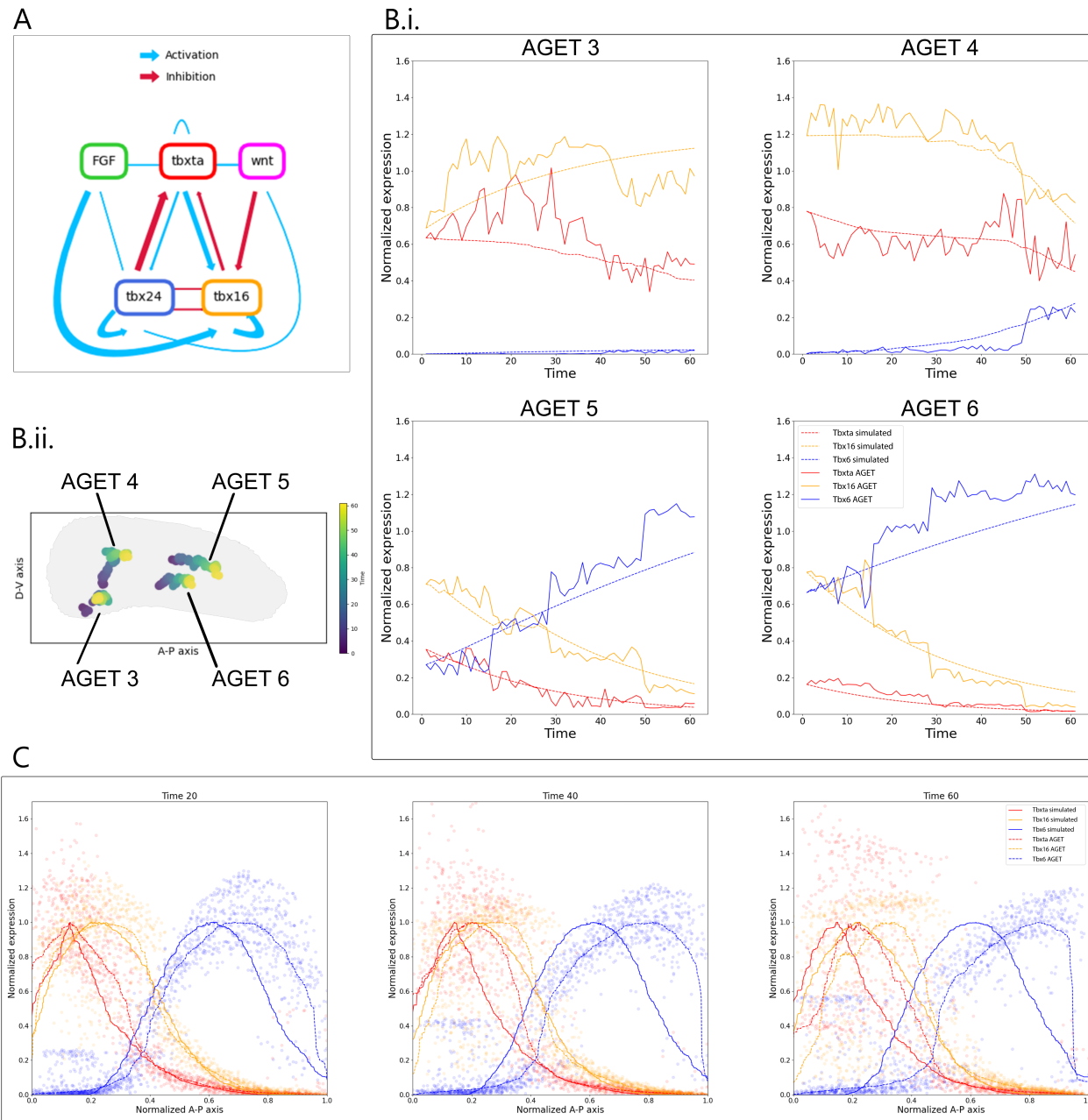


Fig. 3. Performance and fit of the GRN corresponding to the maximum a posteriori (MAP) parameters. (A) GRN topology with MAP parameters obtained from the MCMC inference. **(B)i.** Simulated data (dotted curves) for four of the ten AGETs (solid curves) used for model fitting **(B)ii.** Illustrative spatial location in the PSM of the four AGETs shown in B.i. **(C)** Snapshots showing simulated T-box gene expression along the normalized posterior (0) to anterior (1) axis of the PSM at 33%, 66% and 100% of total simulation time respectively. (The full simulation and a quantitative comparison are shown in Supplementary Movies 3 and 4). Dots correspond to the simulated T-box level in a given cell at a given position. The dotted curves have been obtained by fitting smooth curves to the data simulated in all single cells (dots) at each separate time point and normalised. Solid curves have been obtained by fitting smooth curves to the AGETs at each separate time point and normalising in the same way as was done for the simulated data.

aspects of the pattern evolution might be most strongly working on.

While in the previous section we analysed the network corresponding to the parameter set with the maximal posterior probability (MAP) to assess the goodness of fit of one of the candidate GRNs, in this section we assess how well the posterior distribution has been approximated across candidate GRNs (Fig.4). To do this, we selected 1000 parameter sets at random from the posterior

distribution, representing 1000 distinct candidate networks. We then proceeded to cluster them according to the similarity of their parameter values using agglomerative hierarchical clustering (see Materials and Methods). In order to be able to choose a representative to explore further for each cluster, we set the condition that the parameter distributions within clusters should be uni-modal. After imposing this condition, the algorithm returned 30 clusters and the

617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672

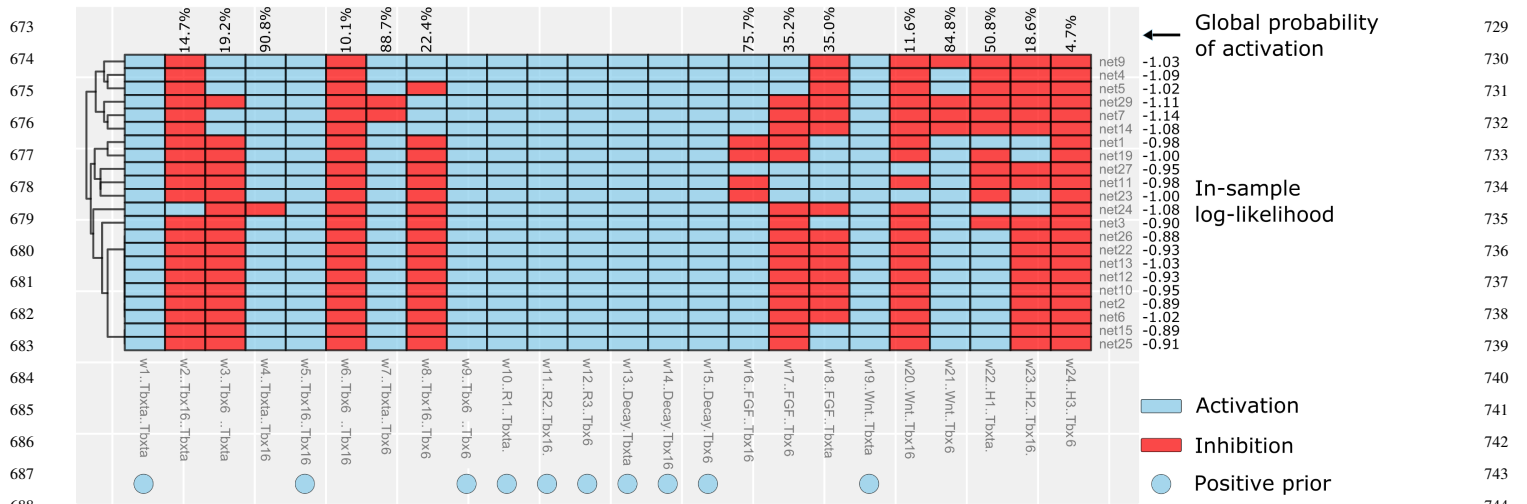


Fig. 4. GRN clusters. The topologies of the mean networks are shown for the 22 well-fitting clusters recovered by the fitting. Rows correspond to representative networks from each cluster, columns represent individual GRN parameters. Quantitative parameters are reduced to whether they are positive or negative for illustration purposes. This can give the impression that some networks and clusters are the same, when in fact they are quantitatively distinct. The percentage above a given parameter indicates the probability that said parameter is positive across clusters. Parameters marked with a blue circle were defined as positive by the prior. In-sample log-likelihood for each network is provided as a measure of goodness of fit.

network with mean parameter values was picked as the representative for each cluster. We simulated the resulting 30 networks and compared them with AGETs 1-10 used for fitting. The simulations were visually inspected and networks returning aberrant patterns were discarded along with all the networks in the cluster that they belonged to. This process left a total of 22 clusters of well-fitting GRNs (Fig.4).

Fig.4 shows the topology of the representative GRNs in each of the resulting 22 clusters. By topology we mean whether parameters are positive (blue) or negative (negative). This provides only a superficial illustration of the clusters which, while useful for visualisation purposes, omits much of the complexity within these classes since the clustering was done on the quantitative value of the parameters. For this reason too, it might appear that representative networks of different clusters are the same, however although that might be the case qualitatively (taking only into account parameter signs), it isn't the case quantitatively (for example networks 26, 22, 13, 12, 10, 2 and 6). 10 out of 24 parameters were set as positive in the priors (Fig.4, round blue circles; see Materials and Methods for justification), the remaining 14, which correspond to parameters that represent the interaction strengths between T-box genes and from Wnt and FGF to the T-box genes, could adopt positive or negative values. The global probability of an activation (positive parameter) is shown above each corresponding column in Fig.4. Generally, for each parameter there is a clear preference across all clusters, suggesting a degree of constraint in the determinability of parameter values. We also recorded the in-sample log-likelihood of each network as a measure of how well these networks fit the data (Fig.4, right). Given how close these values are, we want to emphasise at this point that they should all be treated as likely candidates and that further biological knowledge and experiments are required to discriminate between them (26). In addition, the number of AGETs used for fitting does not seem to affect the general distribution of parameter values, although it can narrow down the spread of the posterior distributions (Supplementary Figures 5 and 6).

CONCLUSION

Earlier reverse-engineering frameworks have been unable to accommodate the role of cell rearrangements and tissue shape changes in developmental pattern formation. This limitation has heavily biased quantitative studies of pattern formation towards systems where the timing of pattern formation and morphogenesis can be separated. However, the vast majority of patterning processes in animal development do not meet this criterion and in consequence, their study has been grossly under-represented in the GRN literature. As a result, most of our collective knowledge and understanding of the generation and evolution of developmental patterns has been constructed on the omission of any role that might be played by cell movements, tissue shape changes and other morphogenetic mechanisms.

This need not be the case going forward. Thanks to recent advancements in live-imaging and spatial gene expression quantification, the data required to adopt the reverse-engineering framework presented in this paper is becoming available in an ever-increasing number of species spanning the range of animal phylogeny. This will make it possible to construct AGETs and infer GRNs in a wider range of systems. Simulation and subsequent analysis of patterning processes that are dependent on or at least, co-occurring with cell movements will increase our understanding of pattern formation and its evolution, and uncover previously hidden general principles that weren't accessible from the restricted types of systems that we were studying. Furthermore, this methodology will find applications well-beyond beyond the study of developmental evolution. In particular, we anticipate a warm reception from fields such as bio-engineering, regenerative medicine and organoid biology, where understanding how 3D cell cultures should be shaped and constrained as they grow to obtain the desired final organisation is paramount and has proven not at all trivial.

Finally, our methodology for the construction of AGETs provides a way in which to visualise approximated gene expression dynamics and patterning without the need for fluorescent transgenic reporter lines, offering an alternative in the form of in-silico

reporters. Once generated, *in silico* reporter lines require no further use of live animals, resulting in a dramatic reduction of the number of animals used in research. In addition, there is in principle no limit to the number of genes that can be reported by an *in silico* reporter line, *in silico* reporter lines could be readily extended to non-model organisms, and they have the potential to exhibit a higher fidelity to the actual dynamics of gene expression since they bypass fluorescent reporter readouts altogether.

MATERIALS AND METHODS

Animal lines and husbandry

This research was regulated under the Animals (Scientific Procedures) Act 1986 Amendment Regulations 2012 following ethical review by the University of Cambridge Animal Welfare and Ethical Review Body (AWERB). Embryos were obtained and raised in standard E3 media at 28°C. Wild Type lines are either Tupfel Long Fin (TL), AB or AB/TL. The Tg(7xTCF-Xla.Sia:GFP) reporter line (33) was provided by the Steven Wilson laboratory. Embryos were staged as in (34).

In Situ Hybridisation Chain Reaction (HCR)

Embryos were incubated until they reached the the desired developmental stage, then fixed in 4% PFA in DEPC treated PBS without calcium and magnesium, and stored at 4°C overnight. Once fixed, embryos were stained using HCR version 3 following the standard zebrafish protocol found in (29). Probes, fluorescent hairpins and buffers were all purchased from Molecular Instruments. After staining, samples were stained with DAPI and mounted using 80% glycerol.

Immunohistochemistry

Embryos were incubated until they reached the desired developmental stage, then fixed in 4% PFA in DEPC treated PBS without calcium and magnesium, and stored at 4°C overnight. The embryos were subsequently blocked in 3% goat serum in 0.25% Triton, 1% DMSO, in PBS for one hour at room temperature. Our read-out for FGF activity - Diphosphorylated ERK - was detected using the primary antibody (M9692-200UL, Sigma) diluted 1 in 500 in 3% goat serum in 0.25% Triton, 1% DMSO, in PBS. All samples were incubated at 4°C overnight and then washed in 0.25% Triton, 1% DMSO, in PBS. Secondary Alexa 647nm conjugated antibodies were diluted 1 in 500 in 3% goat serum in 0.25% Triton, 1% DMSO, 1X DAPI in PBS and applied overnight at 4°C.

Imaging and image analysis

Fixed HCR and immunostained samples were imaged with a Zeiss LSM700 inverted confocal at 12 bit, using either 20X or 40X magnification and an image resolution of 512x512 pixels. Nuclear segmentation of whole stained embryonic tailbuds was performed using a tight mask applied around the DAPI stain using Imaris (Bitplane) with a surface detail of 0.5µm. Positional values for each nucleus were exported as X, Y, Z coordinates relative to the posterior-most tip of the PSM where X, Y, Z were equal to (0, 0, 0). The PSM was then segmented by hand by deleting nuclear surfaces outside of the PSM, including notochord, spinal cord, anterior somites and ectoderm. PSM length was normalised individually between 0 and 1 by division of the position in X by the maximum X value measured in each embryo.

Single cell image analysis was conducted using Imaris (Bitplane) by generating loose surface masks around the DAPI stain to capture the full nuclear region and a small region of cytoplasm. Surface masks were then filtered to remove any masks where two cells joined together or small surfaces caused by background noise, or fragmented apoptotic nuclei. The intensity sum of each channel was measured and normalised by the area of the surface. Expression level was then normalised between 0 and 1 using the maximum value measured for each gene, in each experiment.

Live imaging datasets of the developing PSM were created using a TriM Scope II Upright 2-photon scanning fluorescence microscope equipped Insight DeepSee dual-line laser (tunable 710-1300 nm fixed 1040 nm line) (see details in (27)). The developing embryo was imaged with a 25X 1.05 NA water dipping objective. Embryos were positioned laterally in low melting agarose with the entire tail cut free to allow for normal development (35). Tracks were generated automatically and validated manually using the Imaris imaging software.

Aligning point clouds with ICP

We used the Python library Open3d (36) and the implementation of the point-to-plane ICP (Iterative Closest Point) algorithm therein (30) to perform the point cloud alignment. ICP algorithms can be used to align two point clouds from an initial approximate alignment. The aim is to find a transformation matrix, that rotates and moves the source point cloud in a way that achieves an optimal alignment with the target point cloud. ICP algorithms work by iterating two steps. First, for each point in the source point cloud, the algorithm will determine the corresponding closest point in the target point cloud. Second, the algorithm will find the transformation matrix that most optimally minimizes the distances between the corresponding points. The result is a transformed source point cloud that is closely aligned with the target point cloud. As a pre-processing step, the source and target point clouds have been re-scaled to have the same A-P length. Since we are working with biological tissues, point clouds will not correspond exactly, differing slightly in size and shape. This will impact the quality of the resulting alignment which had to be visually assessed and validated. In this case study, three of the thirteen source images were excluded from the analysis due to poor alignment.

AGET construction

While the main methodology used for constructing AGETs is covered in the results section, below (Algorithm 1) we provide pseudo-code that describes the same process.

Mathematical model formulation

We used a dynamical systems formulation model the T-box gene regulatory network in the zebrafish PSM. The model's aim is to recapitulate the dynamics of T-box gene expression in every cell in the developing zebrafish PSM, generating the emergence of the tissue-level T-box gene expression pattern. We use a connectionist model formulation which has been extensively used and validated to previously model other developmental patterning processes (37; 14; 8).

The mRNA concentrations encoded by the T-box genes *tbxta*, *tbx16* and *tbx6* are represented by the state variables of the dynamical system. For each gene, the concentration of its associated mRNA a at time t is given by $g^a(t)$. mRNA concentration over

Algorithm 1: Mapping T-box gene expression from HCR images onto tracking data

Result: AGETs: Cell tracks with dynamic T-box and signalling expression information

Create target point clouds from tracking data $Target_i$, for every time point $i \in 1, \dots, 61$;

Create source point clouds with gene expression information from HCR data $Source_j$, for every source image $j \in 1, \dots, 10$;

for i in $1 : 61$ **do**

for j in $1 : 10$ **do**

 Align $Source_j$ and $Target_i$ using ICP registration;

for $Cell_k$ in $Target_i$ **do**

 Find $n=5$ closest neighbours of $Cell_k$ in $Source_j$;

 Calculate median M_{ijk} of closest neighbours;

 Assign M_{ijk} to $Cell_k$;

end

end

for $Cell_k$ in $Target_i$ **do**

 Calculate median M_{ik} of medians M_{ijk} from 10 source point clouds $Source_{1:10}$;

 Assign M_{ik} to $Cell_k$;

end

end

Extract all cell tracks with their assigned gene expression (AGETs)

time is governed by the following system of three coupled ordinary differential equations:

$$\frac{dg_a(t)}{dt} = R_a \phi(u_a) - \lambda_a g_a(t) \quad (1)$$

where R^a and λ^a respectively represent the rates of mRNA production and decay. ϕ is a sigmoid regulation-expression function used to represent the cooperative, saturating, coarse-grained kinetics of transcriptional regulation and introduces non-linearities into the model that enable it to exhibit complex dynamics:

$$\phi(u_a) = \frac{1}{2} \left(\frac{u_a}{\sqrt{(u_a)^2 + 1}} + 1 \right), \quad (2)$$

where

$$u_a = \sum_{b \in G} W^{ba} g_b(t) + \sum_{s \in S} E^{sa} g_s(t) + h_a. \quad (3)$$

$G = \{tbxta, tbx16, tbx6\}$ refers to the set of T-box genes while $S = \{\text{Wnt}, \text{FGF}\}$ represents the set of external regulatory inputs provided by the Wnt and FGF signalling environments. The concentrations of the external regulators g_s are provided directly from the AGETs into the simulation and are not themselves being modelled. Changing Wnt and FGF concentrations over time renders the parameter term $\sum_{s \in S} E^{sa} g_s(t)$ time-dependent and therefore, the model non-autonomous (38; 39).

The inter-connectivity matrices W and E house the parameters representing the regulatory interactions among the T-box genes, and from Wnt and FGF to the T-box genes, respectively. Matrix

elements w^{ba} and e^{sa} are the parameters representing the effect of regulator b or s on target gene a . These can be positive (representing an activation from b or s onto a), negative (representing a repression), or close to zero (no interaction). h_a is a threshold parameter denoting the basal activity of gene a , which acknowledges the possible presence of regulators absent from our model. To perform the live-modelling simulations, the same model formulation is implemented in each cell in the time-lapse. Initial concentrations of $tbxta$, $tbx16$ and $tbx6$ are read out directly from the first time point of the AGET corresponding to that cell, and dynamic Wnt and FGF values are updated from the same AGET.

Model fitting: MCMC approach

We used the Markov Chain Monte Carlo approach implemented in the Python emcee library (40) to approximate the posterior distribution of the GRN parameters. A property of this implementation is the use of an ensemble of walkers, rather than a single one. To fit to 10 AGETs, we used a uniform prior from -200 to +200, except when the prior were restricted, and fitted to the time scale used in the simulation. The time scale was chosen such that 1 equals the time that the fastest cell takes to travel through the whole PSM and enter a somite. We used a Gaussian distribution with fixed standard deviations per gene to model the differences between simulated gene expression and target gene expression approximated by the AGETs, and in this way obtain a likelihood function. We ran the MCMC with 70 walkers and for a total of 50'000 steps. Although the auto-correlation time was high and the acceptance fraction with 4.1% was on the low side, the inferred parameters led to well-fitting simulated data. Model training took approximately three days using 20 cores. To generate the supplementary information figures where we assess the performance of fitting to different numbers of AGETs, we use the same range for the prior distributions, but this time with 100 walkers and 10,000 steps. Supplementary Figure 7 shows how the mean acceptance fraction increases per run with the number of AGETs used and the mean auto-correlation score per run decreases as the number of AGETs increases until 200 AGETs, and stabilises thereafter.

Acknowledgements

The authors would like to thank the Cambridge Advanced Imaging Centre (CAIC) for imaging support and the University of Oxford Advanced Research Computing (ARC) facility for their computational support.

Competing interests

The authors declare no competing interests.

Contribution

Conceptualisation: BP, BS and BV. Methodology: KS, BV, BS and BP. Software: KS, SET, KT, DS, SH, YW, BP and BV. Validation: KS and SET. Formal Analysis: KS and SET. Investigation (experimental work): TF. Writing-original draft preparation: KS and BV. Writing-review and editing: BS and BV. Supervision: BP, BS and BV. Funding acquisition: BP, BS and BV.

Funding

K.S. was initially supported by Wave 1 of the UKRI Strategic Priorities Fund under the EPSRC grant EP/T001569/1, particularly the "AI for Science and Government" theme within that grant and the Alan Turing Institute, and later by a by a Henry Dale Fellowship granted to B.S. jointly funded by the Wellcome Trust and the Royal Society (109408/Z/15/Z). S.E.T. was supported by a Clarendon Scholarship. T.F., S.H. and B.S. are supported by a Henry Dale Fellowship jointly funded by the Wellcome Trust and the Royal Society (109408/Z/15/Z) and T.F. by a scholarship from the Cambridge Trust, University of Cambridge. L.T. is supported by a scholarship from the BBSRC. Y.W. is supported by a summer vacation stipend from St Catharine's College, University of Cambridge. B.C. is supported by a Stipend from the Bedford Fund, King's College, University of Cambridge and a scholarship from the Wellcome Trust. B.V. was supported by a Herschel Smith Postdoctoral Fellowship, University of Cambridge and Department of Zoology, University of Oxford. B.C. is supported by a Wellcome Trust Developmental Mechanisms PhD studentship (222279/Z/20/Z).

Data availability

Data and code available from <https://github.com/spikay/AGETS>

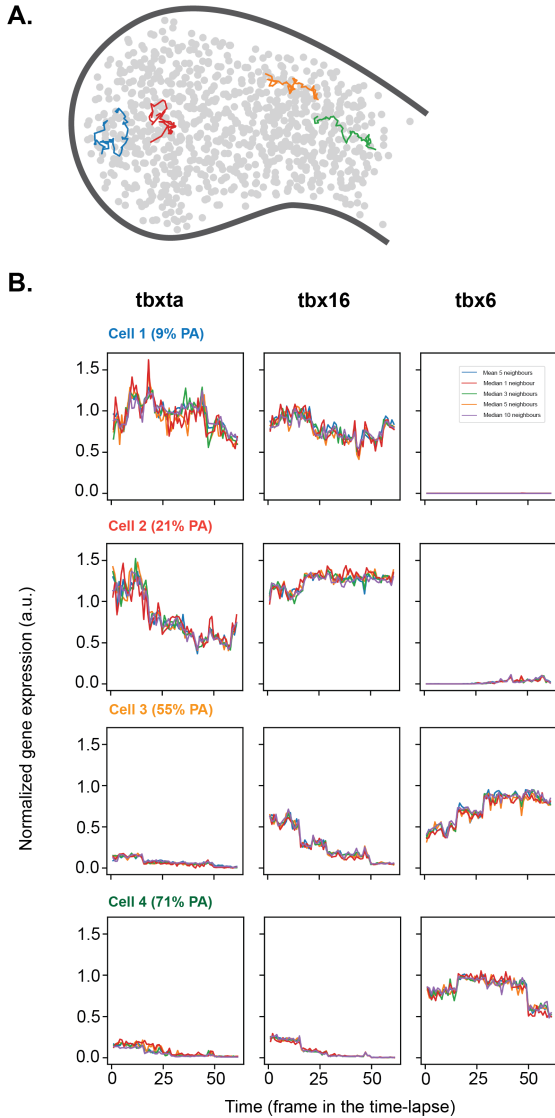
REFERENCES

- [1]Reinitz J, Sharp DH. Gene circuits and their uses. *Integrative Approaches to Molecular Biology*. 1996:253-72.
- [2]Liang S, Fuhrman S, Somogyi R, et al. Reveal, a general reverse engineering algorithm for inference of genetic network architectures. In: *Pacific symposium on biocomputing*. vol. 3. Citeseer; 1998. p. 18-29.
- [3]D'haeseleer P, Liang S, Somogyi R. Genetic network inference: from co-expression clustering to reverse engineering. *Bioinformatics*. 2000;16(8):707-26.
- [4]Gardner TS, Faith JJ. Reverse-engineering transcription control networks. *Physics of life reviews*. 2005;2(1):65-88.
- [5]Rockman MV. Reverse engineering the genotype-phenotype map with natural genetic variation. *Nature*. 2008;456(7223):738-44.
- [6]He F, Balling R, Zeng AP. Reverse engineering and verification of gene networks: principles, assumptions, and limitations of present methods and future perspectives. *Journal of biotechnology*. 2009;144(3):190-203.
- [7]Jaeger J, Monk NA. Reverse Engineering of Gene Regulatory Networks. *Learning and inference in computational systems biology*. 2010;9:34.
- [8]Crombach A, Wotton KR, Cicin-Sain D, Ashyraliyev M, Jaeger J. Efficient reverse-engineering of a developmental gene regulatory network. *PLoS computational biology*. 2012;8(7):e1002589.
- [9]Verd B, Crombach A, Jaeger J. Dynamic maternal gradients control timing and shift-rates for *Drosophila* gap gene expression. *PLOS Computational Biology*. 2017;13(2):e1005285.
- [10]Verd B, Clark E, Wotton KR, Janssens H, Jiménez-Guri E, Crombach A, et al. A damped oscillator imposes temporal order on posterior gap gene expression in *Drosophila*. *PLoS biology*. 2018;16(2):e2003174.
- [11]Manu S, Surkova, Spirov AV, Gursky VV, Janssens H, Kim AR, Radulescu O, et al. Canalization of gene expression and domain shifts in

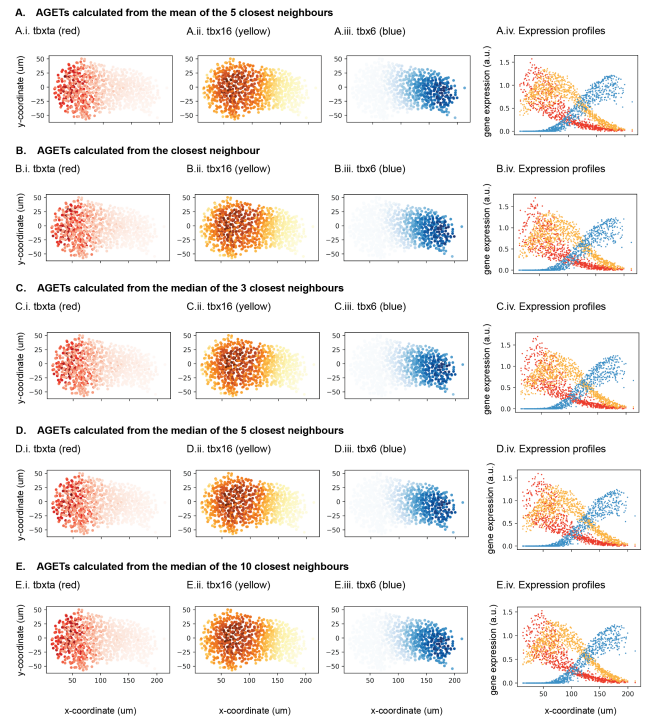
- the *Drosophila* blastoderm by dynamical attractors. *PLoS computational biology*. 2009;5(3):e1000303.
- [12]Balaskas N, Ribeiro A, Panovska J, Dessaud E, Sasai N, Page KM, et al. Gene regulatory logic for reading the Sonic Hedgehog signaling gradient in the vertebrate neural tube. *Cell*. 2012;148(1-2):273-84.
- [13]Verd B, Monk NA, Jaeger J. Modularity, criticality, and evolvability of a developmental gene regulatory network. *Elife*. 2019;8:e42832.
- [14]Jaeger J, Blagov M, Kosman D, Kozlov KN, Manu, Myasnikova E, et al. Dynamical analysis of regulatory interactions in the gap gene system of *Drosophila melanogaster*. *Genetics*. 2004;167(4):1721-37.
- [15]Cohen M, Kicheva A, Ribeiro A, Blassberg R, Page KM, Barnes CP, et al. Ptc1 and Gli regulate Shh signalling dynamics via multiple mechanisms. *Nature communications*. 2015;6(1):1-12.
- [16]Kicheva A, Bollenbach T, Ribeiro A, Valle HP, Lovell-Badge R, Episkopou V, et al. Coordination of progenitor specification and growth in mouse and chick spinal cord. *Science*. 2014;345(6204).
- [17]El-Sherif E, Zhu X, Fu J, Brown SJ. Caudal regulates the spatiotemporal dynamics of pair-rule waves in *Tribolium*. *PLoS genetics*. 2014;10(10):e1004677.
- [18]Wu H, Jiao R, Ma J, et al. Temporal and spatial dynamics of scaling-specific features of a gene regulatory network in *Drosophila*. *Nature communications*. 2015;6(1):1-13.
- [19]Averbukh I, Lai SL, Doe CQ, Barkai N. A repressor-decay timer for robust temporal patterning in embryonic *Drosophila* neuroblast lineages. *Elife*. 2018;7:e38631.
- [20]Schröter C, Ares S, Morelli LG, Isakova A, Hens K, Soroldoni D, et al. Topology and dynamics of the zebrafish segmentation clock core circuit. *PLoS biology*. 2012;10(7):e1001364.
- [21]Rayon T, Stamatakis D, Perez-Carrasco R, Garcia-Perez L, Barrington C, Melchionda M, et al. Species-specific pace of development is associated with differences in protein stability. *Science*. 2020;369(6510).
- [22]Uriu K, Morishita Y, Iwasa Y. Random cell movement promotes synchronization of the segmentation clock. *Proceedings of the National Academy of Sciences*. 2010;107(11):4979-84.
- [23]Crombach A, Wotton KR, Jiménez-Guri E, Jaeger J. Gap gene regulatory dynamics evolve along a genotype network. *Molecular biology and evolution*. 2016;33(5):1293-307.
- [24]Kicheva A, Cohen M, Briscoe J. Developmental pattern formation: insights from physics and biology. *Science*. 2012;338(6104):210-2.
- [25]Huch M, Knoblich JA, Lutolf MP, Martínez-Arias A. The hope and the hype of organoid research. *Development*. 2017;144(6):938-41.
- [26]Fulton T, Speiss K, Thomson L, Wang Y, Clark B, Hwang S, et al. Cell Rearrangement Generates Pattern Emergence as a Function of Temporal Morphogen Exposure. *bioRxiv*. 2022.
- [27]Thomson L, Muresan L, Steventon B. The zebrafish presomitic mesoderm elongates through compaction-extension. *Cells & Development*. 2021.
- [28]Mongera A, Rowghanian P, Gustafson HJ, Shelton E, Kealhofer DA, Carn EK, et al. A fluid-to-solid jamming transition underlies vertebrate body axis elongation. *Nature*. 2018;561(7723):401-5.
- [29]Choi HM, Schwarzkopf M, Fornace ME, Acharya A, Artavanis G, Stegmaier J, et al. Third-generation in situ hybridization chain reaction: multiplexed, quantitative, sensitive, versatile, robust. *Development*. 2018;145(12):dev165753.
- [30]Rusinkiewicz S, Levoy M. Efficient variants of the ICP algorithm. In: *Proceedings third international conference on 3-D digital imaging and modeling*. IEEE; 2001. p. 145-52.
- [31]Ram R, Chetty M. MCMC based Bayesian inference for modeling gene networks. In: *IAPR International Conference on Pattern Recognition in Bioinformatics*. Springer; 2009. p. 293-306.
- [32]Kawamura A, Koshida S, Hijikata H, Ohbayashi A, Kondoh H, Takada S. Groucho-associated transcriptional repressor ripply1 is required for proper transition from the presomitic mesoderm to somites. *Developmental cell*. 2005;9(6):735-44.
- [33]Moro E, Ozhan-Kizil G, Mongera A, Beis D, Wierzbicki C, Young RM, et al. In vivo Wnt signaling tracing through a transgenic biosensor fish reveals novel activity domains. *Developmental biology*. 2012;366(2):327-40.

1121	[34]Kimmel CB, Ballard WW, Kimmel SR, Ullmann B, Schilling TF. Stages of embryonic development of the zebrafish. <i>Developmental dynamics</i> . 1995;203(3):253-310.	1177
1122		1178
1123	[35]Hirsinger E, Steventon B. A versatile mounting method for long term imaging of zebrafish development. <i>JoVE (Journal of Visualized Experiments)</i> . 2017;(119):e55210.	1179
1124		1180
1125		1181
1126	[36]Zhou QY, Park J, Koltun V. Open3D: A Modern Library for 3D Data Processing. arXiv:180109847. 2018.	1182
1127		1183
1128	[37]Mjolsness E, Sharp DH, Reinitz J. A connectionist model of development. <i>Journal of theoretical Biology</i> . 1991;152(4):429-53.	1184
1129		1185
1130	[38]Collier JR, Monk NA, Maini PK, Lewis JH. Pattern formation by lateral inhibition with feedback: a mathematical model of delta-notch intercellular signalling. <i>Journal of theoretical Biology</i> . 1996;183(4):429-46.	1186
1131		1187
1132	[39]Verd B, Crombach A, Jaeger J. Classification of transient behaviours in a time-dependent toggle switch model. <i>BMC systems biology</i> . 2014;8(1):1-19.	1188
1133		1189
1134	[40]Foreman-Mackey D, Hogg DW, Lang D, Goodman J. emcee: the MCMC hammer. <i>Publications of the Astronomical Society of the Pacific</i> . 2013;125(925):306.	1190
1135		1191
1136		1192
1137		1193
1138		1194
1139		1195
1140		1196
1141		1197
1142		1198
1143		1199
1144		1200
1145		1201
1146		1202
1147		1203
1148		1204
1149		1205
1150		1206
1151		1207
1152		1208
1153		1209
1154		1210
1155		1211
1156		1212
1157		1213
1158		1214
1159		1215
1160		1216
1161		1217
1162		1218
1163		1219
1164		1220
1165		1221
1166		1222
1167		1223
1168		1224
1169		1225
1170		1226
1171		1227
1172		1228
1173		1229
1174		1230
1175		1231
1176		1232

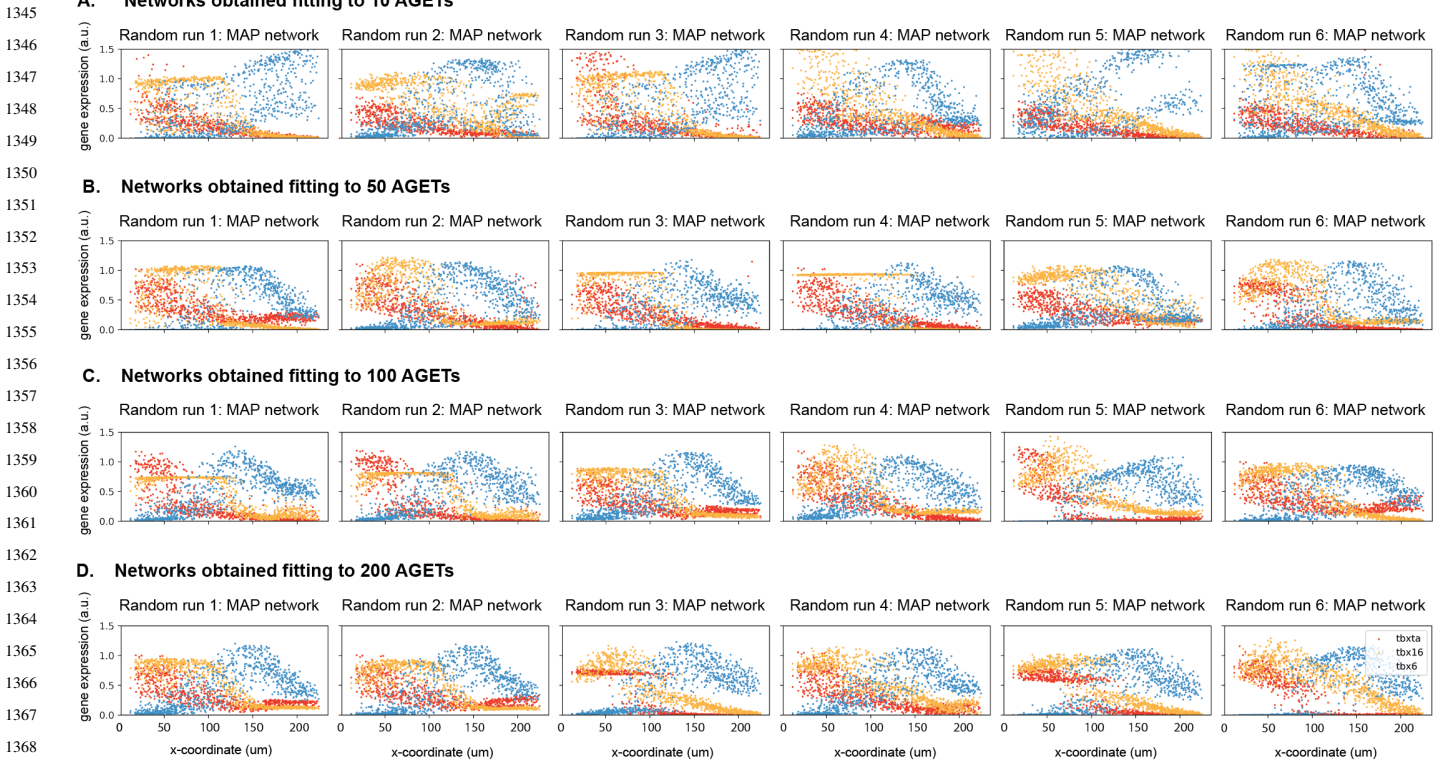
Supplementary



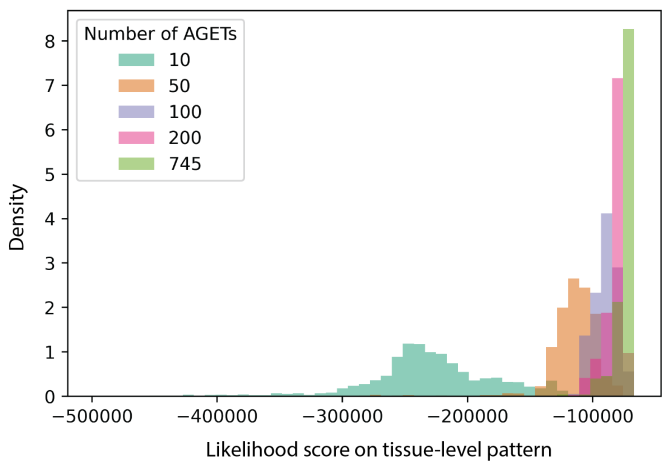
Supplementary Fig. 1. AGETs are robust to calculation method. (A) The spatial trajectories of four cells used in (B) to test AGET robustness to changing the rules used to calculate them. (B) Each row shows the AGET values calculated for a cell (cell 1; located initially at 9% PA position, cell 2; located initially at 21% PA position, cell 3; located initially at 55% PA position and cell 4, located initially at 71% PA position) for tbxta (column 1), tbx16 (column 2) and tbx6 (column 3). AGETs are calculated taking the mean of the five nearest neighbours (blue), the median of the nearest neighbour (red), the median of the three nearest neighbours (green), the median of the five nearest neighbours (orange) or the median of the ten nearest neighbours (purple).



Supplementary Fig. 2. Tissue-level pattern is robust to AGET calculation method. (A) Approximated Tbox gene expression pattern on the PSM when AGETs were calculated taking the mean of the five nearest neighbours. (A.i.) Approximated tbxta in the cells of the PSM. Each dot represents the position of a cell in the PSM (x,y-projection where dorsal is to the top and posterior is to the left). Shade of red indicates approximated tbxta concentration (dark red, highest, white, lowest). (A.ii.) Approximated tbx16 in the cells of the PSM. Shade of yellow indicates approximated tbx16 concentration (dark yellow, highest, white, lowest). (A.iii.) Approximated tbx6 in the cells of the PSM. Shade of blue indicates approximated tbx6 concentration (dark blue, highest, white, lowest). (A.iv.) Tbox gene expression profiles. Each dot represents the concentration of one of the tbox genes (tbxta (red), tbx16 (yellow) and tbx6 (blue)) in a given cell. The position along the posterior to anterior axis of each cell is given by its x-coordinate. (B) Approximated Tbox gene expression pattern on the PSM when AGETs were calculated taking the value of the nearest neighbour. (B.i.) Approximated tbxta in the cells of the PSM. (B.ii.) Approximated tbx16 in the cells of the PSM. (B.iii.) Approximated tbx6 in the cells of the PSM. (B.iv.) Tbox gene expression profiles. (C) Approximated Tbox gene expression pattern on the PSM when AGETs were calculated taking the value of the nearest neighbour. (C.i.) Approximated tbxta in the cells of the PSM. (C.ii.) Approximated tbx16 in the cells of the PSM. (C.iii.) Approximated tbx6 in the cells of the PSM. (C.iv.) Tbox gene expression profiles. (D) Approximated Tbox gene expression pattern on the PSM when AGETs were calculated taking the value of the nearest neighbour. (D.i.) Approximated tbxta in the cells of the PSM. (D.ii.) Approximated tbx16 in the cells of the PSM. (D.iii.) Approximated tbx6 in the cells of the PSM. (D.iv.) Tbox gene expression profiles. (E) Approximated Tbox gene expression pattern on the PSM when AGETs were calculated taking the value of the nearest neighbour. (E.i.) Approximated tbxta in the cells of the PSM. (E.ii.) Approximated tbx16 in the cells of the PSM. (E.iii.) Approximated tbx6 in the cells of the PSM. (E.iv.) Tbox gene expression profiles.

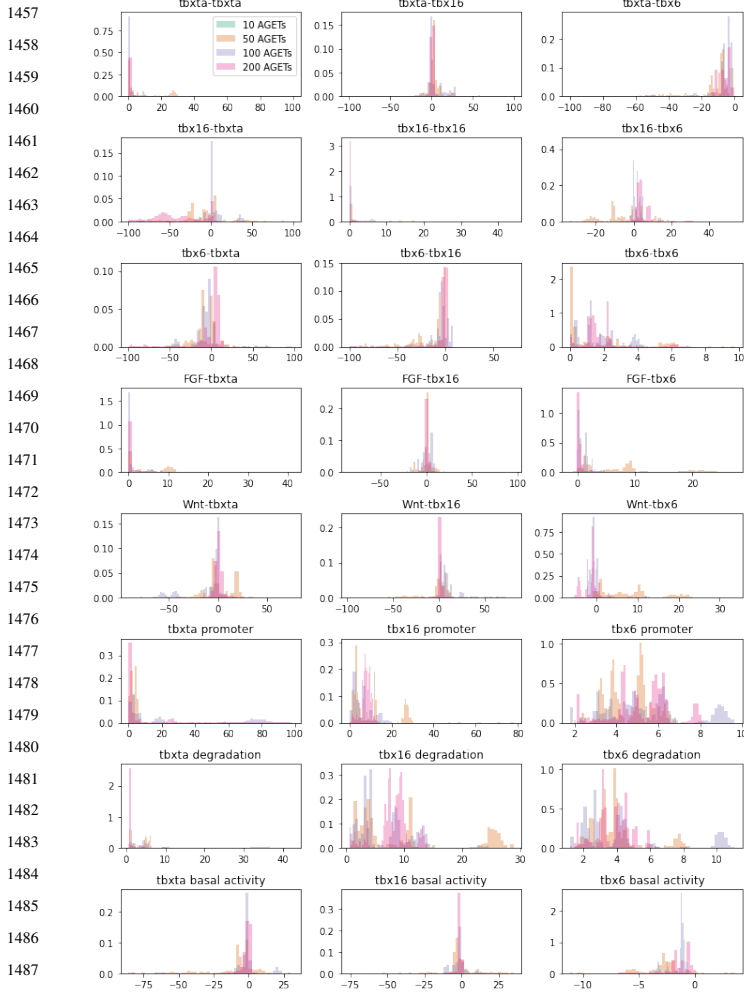


Supplementary Fig. 3. The proportion of parameter combinations producing good fits increases as the number of AGETs used for fitting is increased (A) Networks obtained fitting to 10 AGETs. (B) Networks obtained fitting to 50 AGETs. (C) Networks obtained fitting to 100 AGETs. (D) Networks obtained fitting to 200 AGETs. In each case, the MAP parameter set is taken from 6 independent random runs and the expression profile corresponding to the last time point in the simulation is plotted. Each dot represents the concentration of one of the *tbx* genes (*tbxta* (red), *tbx16* (yellow) and *tbx6* (blue) in a given cell. The position along the posterior to anterior axis of each cell is given by its x-coordinate. Acceptable fits are obtained regardless of the number of AGETs used for fitting, but the proportion of acceptable fits increases with the number of AGETs.

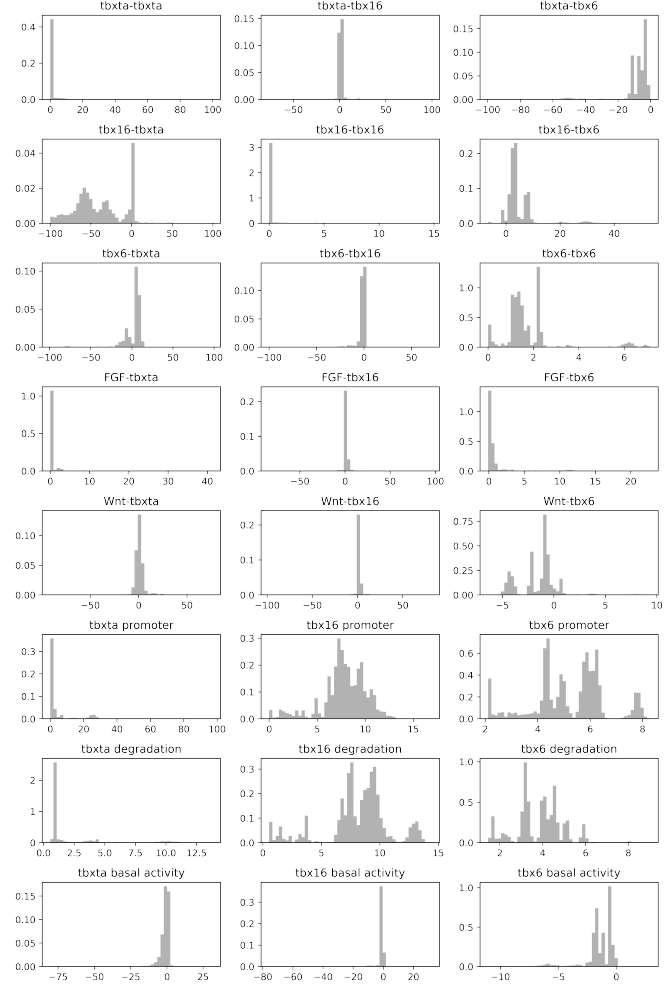


Supplementary Fig. 4. Increasing the number of AGETs used for fitting improves the fits, but good fits are obtained even when fitting to small AGET numbers. Network parameters were inferred using 10, 50, 100, 200 and 745 AGETs to study how AGET number affected the goodness of the fits. For each number of AGETs, 10 rounds of fitting were carried out on three random sets of AGETs, each time using 200 walkers, resulting in 6000 sets of parameter sets inferred per AGET number. For each parameter set, the likelihood score was calculated by comparing the simulated pattern at the level of the tissue with the data. Likelihood scores are colour coded according to the number of AGETs used and plotted on a histogram. Higher likelihood scores reflect better fits. Using more AGETs results in higher average likelihood values and tighter distributions, however high likelihood values are also obtained when smaller numbers of AGETs are used.

1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456



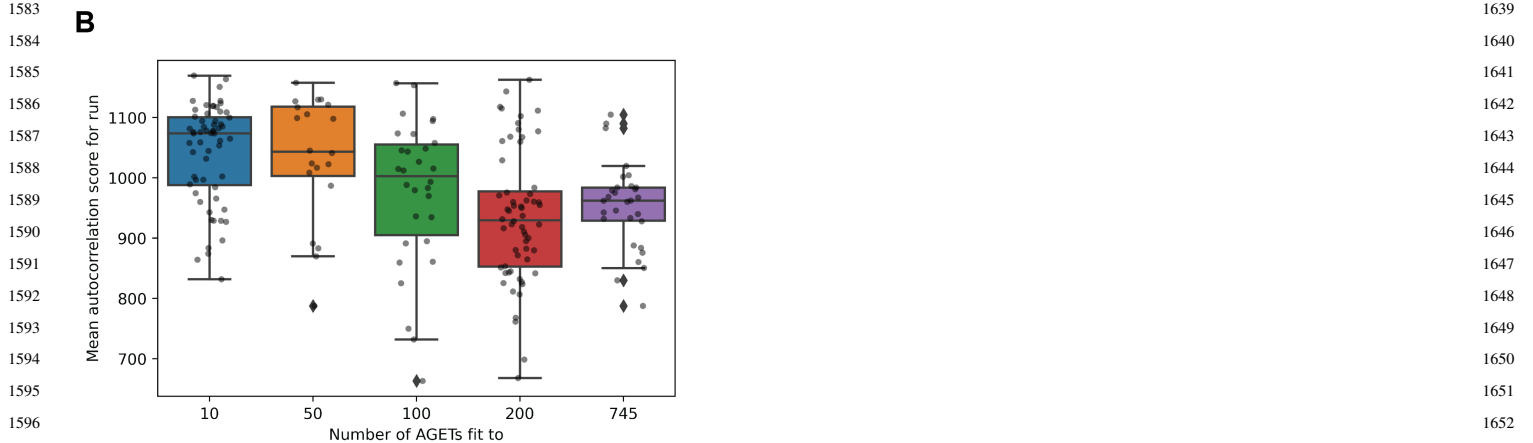
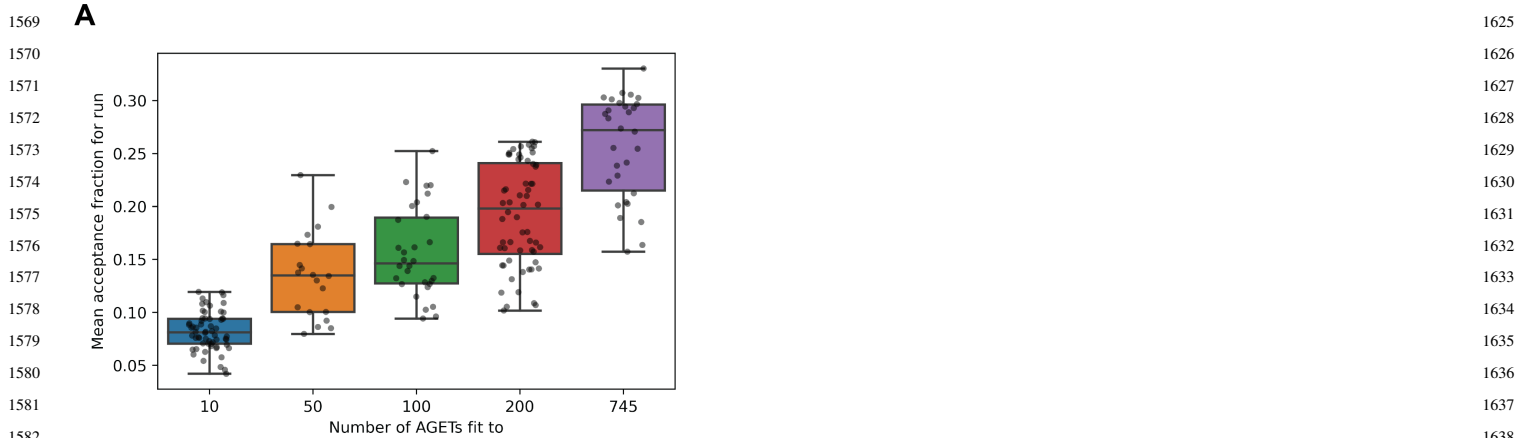
Supplementary Fig. 5. Spread of parameter values obtained using 10, 50, 100 and 200 AGETs.



Supplementary Fig. 6. Spread of parameter values obtained using 200 AGETs.

1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512

1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568



Supplementary Fig. 7. A. Mean acceptance fraction per run increases with the number of AGETs used for fitting. **B.** Mean auto-correlation score per run decreases as the number of AGETs increases until 200 AGETs, and stabilises thereafter.

Supplementary Movie 1. Visualisation of the AGETs for Wnt and FGF on the cell tracks. AGETs were calculated using the median of the five nearest neighbours.

Supplementary Movie 2. Visualisation of the AGETs for Tbxta, Tbx16 and Tbx6 on the cell tracks. AGETs were calculated using the median of the five nearest neighbours.

Supplementary Movie 3. Simulation of the MAP network on the cell tracks.

Supplementary Movie 4. Comparison of simulated and approximated Tbox gene expression on the cell tracks. Coloured dots represent the concentration of tbxta (red), tbx16 (yellow) and tbx6 (blue) in a single cell, plotted against the normalised position of the cell along the PSM, simulated using the MAP network. Dotted lines represent the average simulated tbxta (red), tbx16 (yellow) and tbx6 (blue) domain along the PSM. Solid lines in the bottom right panel represent the average approximated tbxta (red), tbx16 (yellow) and tbx6 (blue) domain along the PSM, where AGETs were calculated using the median of the five nearest neighbours.