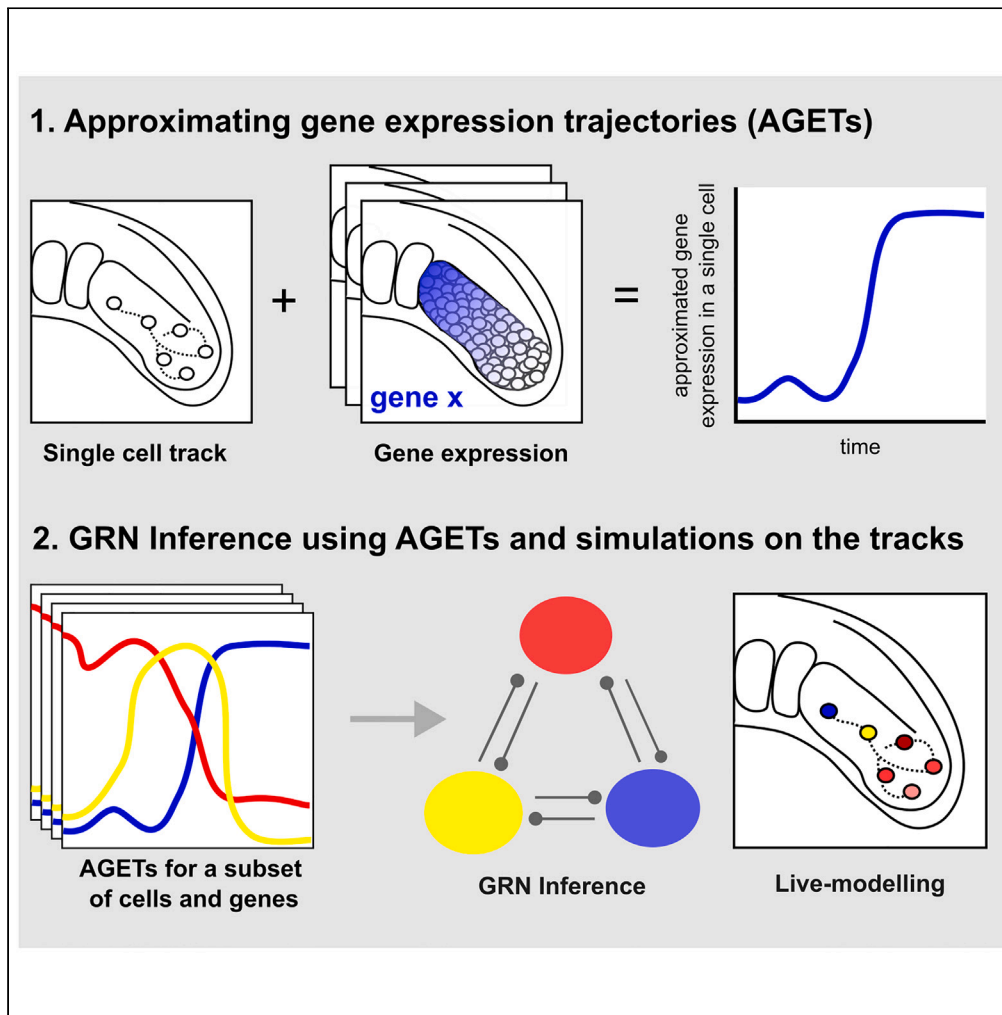


## Article

## Approximated gene expression trajectories for gene regulatory network inference on cell tracks



Kay Spiess,  
Shannon E. Taylor,  
Timothy Fulton, ...,  
Brooks Paige,  
Benjamin  
Steventon, Berta  
Verd

b.paige@ucl.ac.uk (B.P.)  
bjs57@cam.ac.uk (B.S.)  
berta.verdfernandez@biology.  
ox.ac.uk (B.V.)

**Highlights**

We present a method to approximate gene expression trajectories (AGETs) on cell tracks

AGETs are constructed by combining live-imaging data and static gene expression data

AGETs are used to infer the GRNs patterning tissues where the cells are rearranging

GRNs recapitulate patterning, genetic interactions and response to perturbations

## Article

Approximated gene expression trajectories  
for gene regulatory network  
inference on cell tracks

Kay Spiess,<sup>1,2,5</sup> Shannon E. Taylor,<sup>3,5</sup> Timothy Fulton,<sup>1</sup> Kane Toh,<sup>1</sup> Dillan Saunders,<sup>1</sup> Seongwon Hwang,<sup>1</sup>  
Yuxuan Wang,<sup>1</sup> Brooks Paige,<sup>2,4,\*</sup> Benjamin Steventon,<sup>1,\*</sup> and Berta Verd<sup>1,3,6,\*</sup>

## SUMMARY

The study of pattern formation has benefited from our ability to reverse-engineer gene regulatory network (GRN) structure from spatiotemporal quantitative gene expression data. Traditional approaches have focused on systems where the timescales of pattern formation and morphogenesis can be separated. Unfortunately, this is not the case in most animal patterning systems, where pattern formation and morphogenesis are co-occurring and tightly linked. To elucidate patterning mechanisms in such systems we need to adapt our GRN inference methodologies to include cell movements. In this work, we fill this gap by integrating quantitative data from live and fixed embryos to approximate gene expression trajectories (AGETs) in single cells and use these to reverse-engineer GRNs. This framework generates candidate GRNs that recapitulate pattern at the tissue level, gene expression dynamics at the single cell level, recover known genetic interactions and recapitulate experimental perturbations while incorporating cell movements explicitly for the first time.

## 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 spatiotemporal 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.<sup>1–8</sup> Reverse-engineering has been successfully applied to a myriad of systems, from the *Drosophila* blastoderm to the vertebrate neural tube,<sup>9–12</sup> uncovering the mechanisms by which GRNs readout morphogen gradients,<sup>12–17</sup> scale patterns,<sup>18</sup> control the timing of differentiation,<sup>19–21</sup> synchronize cellular fates<sup>22</sup> and evolve pattern formation.<sup>23</sup>

In addition to the spatial patterning of gene expression within tissues, the developing embryo needs to deform and grow tissues to the correct shape and size during morphogenesis. Much of what we know about pattern formation has been learnt from reverse-engineering GRN structure in systems where the timescales of pattern formation and morphogenesis are different and can therefore be separated. In such systems, spatiotemporal gene expression profiles are typically obtained by measuring gene expression levels across the tissue of interest in fixed stained samples, and interpolating between measurements taken at different time points.<sup>8</sup> 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<sup>24</sup> and much insight can be drawn from analyzing reverse-engineered GRNs.<sup>10,13</sup>

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 distributions and proportions are difficult to control as a result of altered patterning due to abnormal morphogenesis in unconstrained tissue geometries.<sup>25</sup> Therefore, in order to be able to understand developmental pattern formation in a broader range of systems, we have to address how morphogenesis and

<sup>1</sup>Department of Genetics, University of Cambridge, Cambridge, UK

<sup>2</sup>The Alan Turing Institute, London, UK

<sup>3</sup>Department of Biology, University of Oxford, Oxford, UK

<sup>4</sup>Centre for Artificial Intelligence, University College London, London, UK

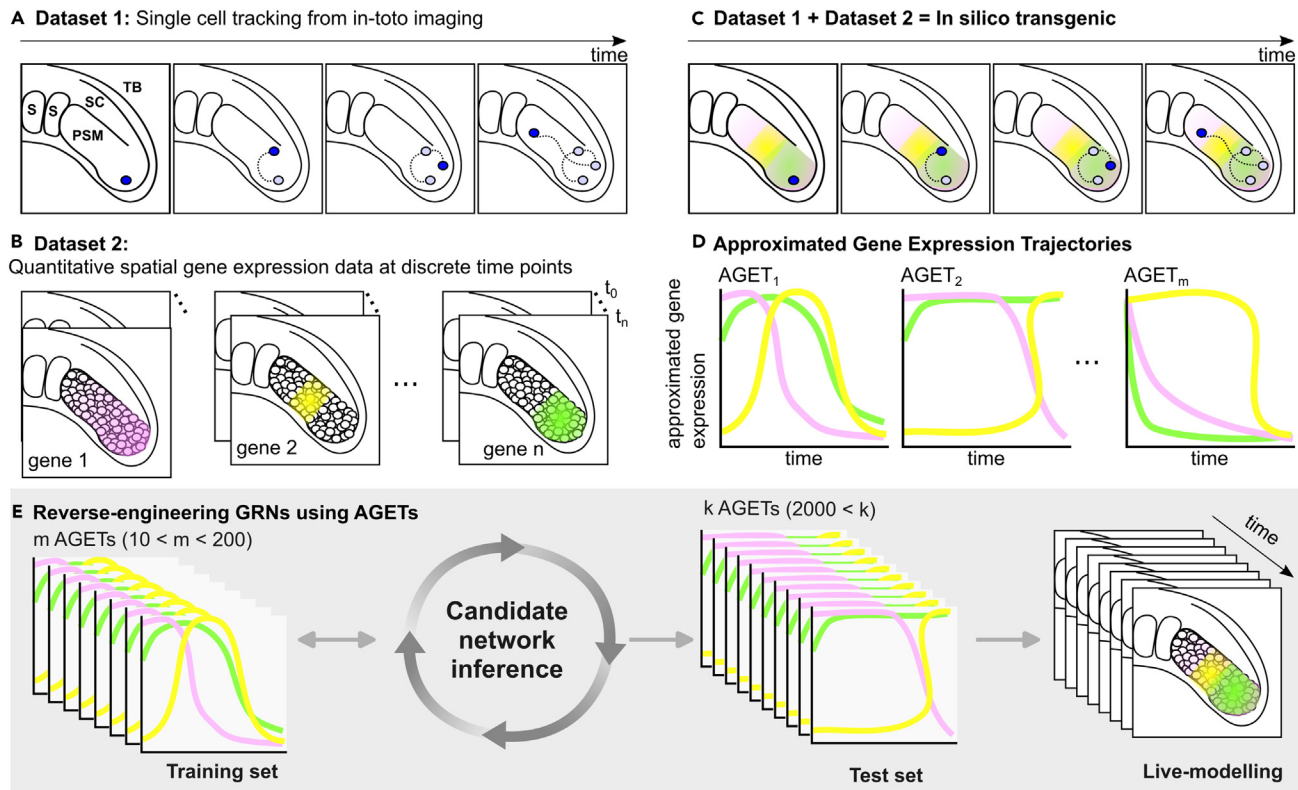
<sup>5</sup>These authors contributed equally

<sup>6</sup>Lead contact

\*Correspondence: [b.paige@ucl.ac.uk](mailto:b.paige@ucl.ac.uk) (B.P.), [bjs57@cam.ac.uk](mailto:bjs57@cam.ac.uk) (B.S.), [berta.verdfernandez@biology.ox.ac.uk](mailto:berta.verdfernandez@biology.ox.ac.uk) (B.V.)

<https://doi.org/10.1016/j.isci.2024.110840>





**Figure 1. Approximated gene expression trajectories (AGETs) to reverse-engineer GRNs patterning tissues undergoing morphogenesis and cell rearrangements**

The generation of AGETs requires the combination of two datasets: (A) single-cell tracking data and (B) HCR gene expression data at the same stages as those covered by the tracking data.

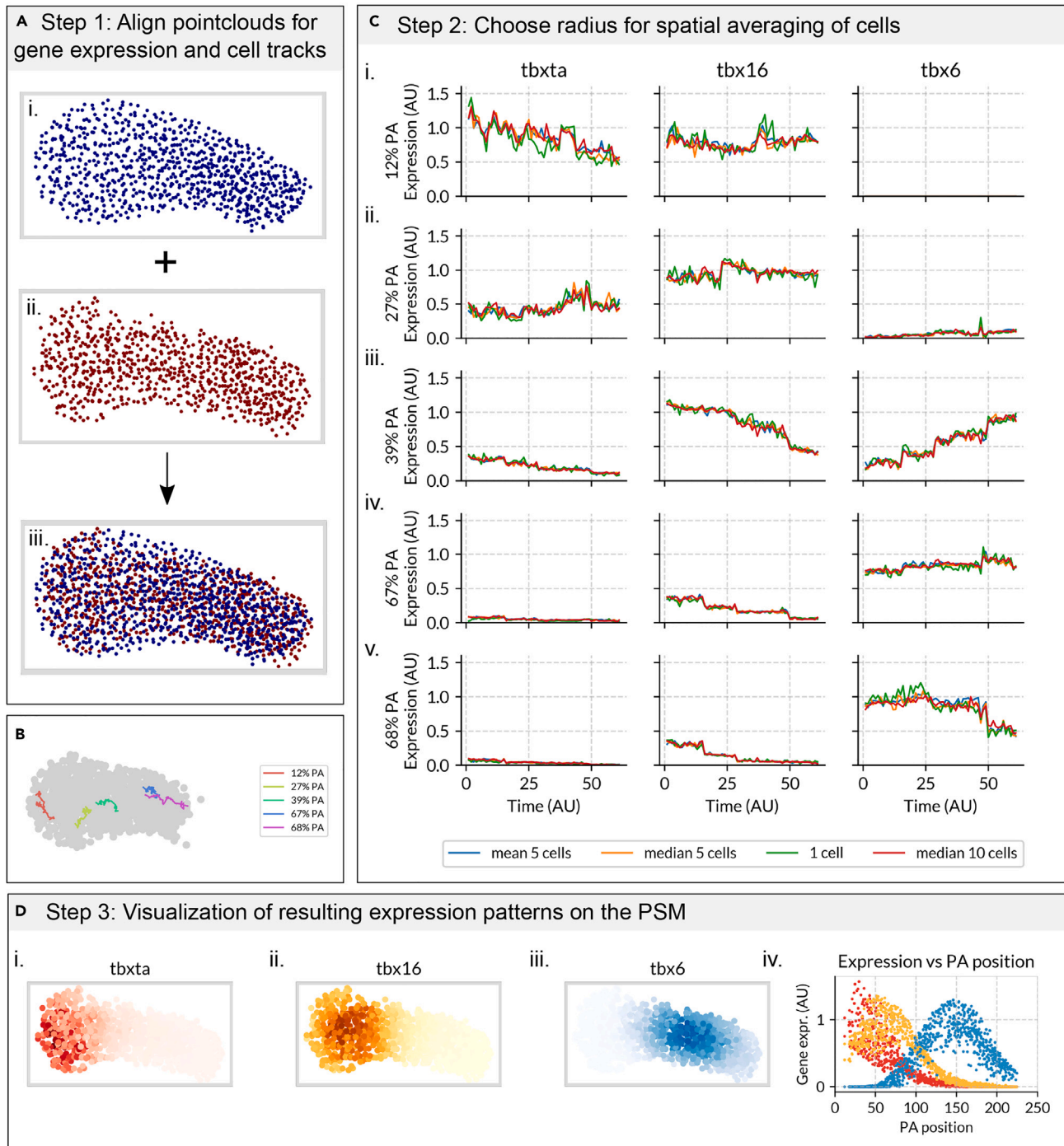
(C) These two datasets are combined to generate an in silico transgenic reporter where every cell in the tracking data has been assigned a gene expression value for every time point using the HCR data.

(E) Subsets of AGETs can then be used to reverse-engineer GRNs, and candidate networks can be simulated using the initial and boundary conditions for all AGETs and represented back on the tracks to recapitulate pattern formation during tissue morphogenesis (live modeling).

GRNs together control fate specification and embryonic organization. Importantly, to be able to do this, we need novel reverse-engineering methodologies that will explicitly accommodate cell movements and tissue shape changes.

Here 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) (see schematic in Figures 1A and S1A). T-box genes coordinate fate specification along the PSM as cells move out of the tailbud and toward the somites.<sup>26</sup> Cell movements in the PSM can be live-imaged and followed in 3D.<sup>27</sup> By the time they reach a somite, cells in the PSM will have undergone a stereotypical progression of T-box gene expression: *tbxta* and *tbx16* in the tailbud, followed by *tbx16* in the posterior PSM and *tbx6* in the anterior PSM (Figures S1A and S1D). 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 behaviors where extensive cell mixing in the tailbud gives way to reduced, almost nonexistent mixing and neighborhood cohesion in the PSM.<sup>28</sup> 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.<sup>26</sup> Despite this, a tissue-level pattern forms which scales with PSM length over the course of posterior axial elongation and somitogenesis.<sup>26</sup> 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, and the cell movements, and tissue shape changes involved in the tissue's morphogenesis.

The reverse-engineering methodology presented in this article 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 (Figure 1A) and three-dimensional quantitative gene expression of the genes and signaling pathways of interest over developmental time (Figure 1B). We project the 3D gene expression data onto the cell tracks to generate Approximated Gene Expression Trajectories (AGETs) in single cells (Figures 1C and 1D). A subset of AGETs is then used to reverse-engineer candidate GRNs applying a Markov Chain Monte Carlo (MCMC) approach (Figure 1E). 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



**Figure 2. AGET Construction**

(A) Step 1: Align point clouds representing the positions of cells' nuclei within the PSM in HCR images and time lapse frames. (A.i) The point cloud in blue represents the positions of the cells (one point, one cell) in the PSM of a processed HCR image (source point cloud). (A.ii) The point cloud in red represents the positions of the cells' nuclei from the first frame of the tracking data (target point cloud). (A.iii) Using ICP, all the source point clouds obtained from the HCR images are aligned with the target point cloud obtained from each frame (61 in total) of the tracking data. This is illustrated by the overlapping red and blue point clouds in the resulting point cloud (bottom).

(B) Relative position within the PSM of the AGETs displayed in C.

(C) Step 2: Assign a gene expression value to every cell in every frame of the tracks based on the expression of its neighbors from the HCRs. AGETs representing approximated T-box gene expression dynamics in cells starting at i. 12%, ii. 27%, iii. 39% iv. 67% and v. 68% PA (posterior-to-anterior) position in the PSM. Y axis

**Figure 2. Continued**

represents relative gene expression levels and x axis, the time frame in the time lapse from 1 to 61. AGETs are calculated using different numbers of neighbors and different averaging methods are shown (blue: a cell is assigned the mean expression of its 5 closest neighbors, orange: a cell is assigned the median expression of its 5 closest neighbors, green: a cell is assigned the expression of its closest neighbor, red: a cell is assigned the median expression of its 10 closest neighbors), and from the AGETs it can be seen that there are no major differences.

(D) Step 3: Visualize the AGETs on the tracks and compare the resulting tissue-level expression patterns with those measured from the HCRs. (D.i) AGETs for *tbxta* (red) D.ii. *tbx16* (yellow) and D.iii. *tbx6* (blue) at the first time frame in the movie. (D.iv) Projection of the the AGETs in D.i-iv. quantified along the posterior to anterior axis of the PSM. Each dot represents the approximated value of a gene (*tbxta* in red, *tbx16* in yellow, and *tbx6* in blue) in a cell. Lines represent maximum projection measures of *tbx* patterns from the HCRs.

expression data, a methodology that we refer to as "live-modelling." The resulting well-fitting GRNs cluster, generating candidate GRNs that we further investigate and challenge using experimental work.<sup>26</sup>

To our knowledge, this inference methodology is the first to integrate cell movements and gene expression data, making it possible to reverse-engineer the 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: approximated gene expression trajectories

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, the current state of the art in live gene expression reporter technology, while very advanced, cannot follow three genes and two signaling 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 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 movie of the developing PSM. The projected expression levels are used to assign gene and signaling 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 recapitulate T-box pattern formation on the developing PSM when simulated on the tracks.

### 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 developing zebrafish tailbud with fluorescently labeled nuclei between the 22nd and 25th somite stages using a two-photon microscope (see<sup>27</sup> and Materials and Methods). Embryonic tailbuds were imaged for 2 h at 2-min intervals, generating 61 consecutive frames. Each frame consists of a point cloud representing the position of single cells in 3D space. All the 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 how individual cells in the PSM move as the zebrafish tailbud develops. 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<sup>29</sup> and antibody stains for the signals Wnt and FGF (see Materials and Methods). As the Tbox pattern scales with tissue growth over developmental time,<sup>26</sup> embryos from different stages could be pooled together, however this method could also be adapted to patterning systems that do not scale with tissue growth. We stained for *tbxta*, *tbx16*, *tbx6*, and DAPI on 23-25ss zebrafish tailbuds using HCR, and were able to quantify 10 of 13 images (2x 23SS, 3x 24SS and 5x 25SS). We measured nuclear gene expression for genes and signals using an Imaris pipeline, again yielding a set of point clouds describing gene expression and nuclear position for each tailbud. See Materials and Methods for a full description of staining, image acquisition, and quantification.

### Approximated gene expression trajectory 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 on 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 onto the tracks to obtain an approximated readout of the gene expression and signaling levels that each cell experiences as it moves.

We first aligned the point clouds representing the positions of the cells in 3D space processed from the HCRs (Figure S1) with the point clouds for each of the 61 time frames in the time lapse (Figure 2A). We use point-to-plane ICP (iterative closest point) to perform this alignment,<sup>30</sup> which in brief, is an iterative algorithm that seeks to map two point clouds onto each other by recursively minimizing the distance

**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)
```

between them (see Materials and Methods). Once the point clouds have been aligned, equivalent regions of different PSMs will overlap in space (Figure 2A) making it possible to use the quantitative gene expression from cells in the processed HCRs to assign gene expression values to the cells in the time lapse at each time frame (Figures 2C and 2D; Algorithm 1 (Materials and Methods)). In this case study, an AGET for a single cell in the tracks will approximate Wnt, FGF, *tbxta*, *tbx16* and *tbx6* values at every time point.

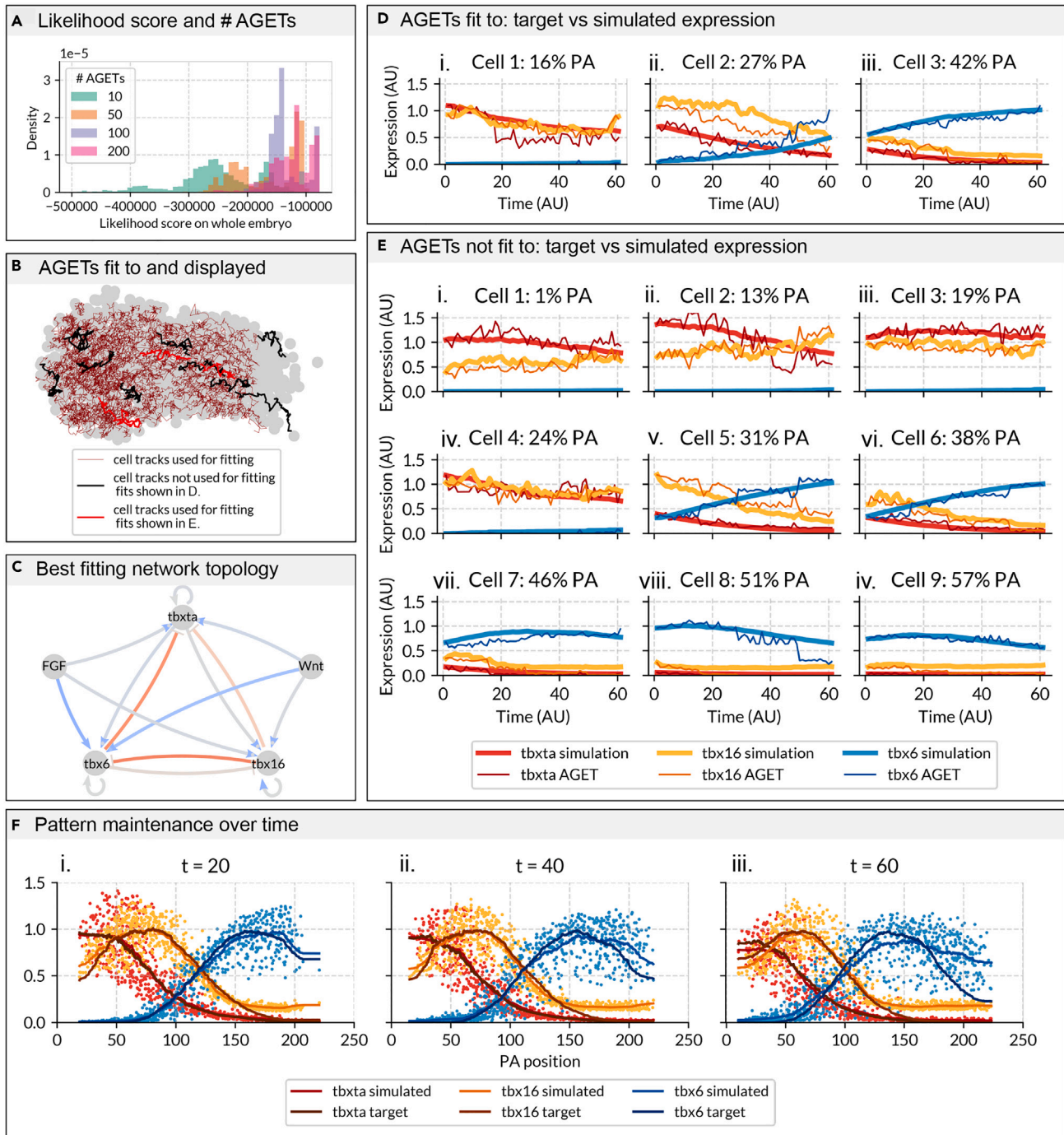
To approximate the gene expression and signaling values in a cell from the time lapse, we first find its  $n$  closest neighboring cells from the processed HCR data (Figures 2B and S2) and assign it an averaged value based on their expression values. This is repeated for every cell in each of the 61 frames in the time lapse, resulting in an AGET for every cell in the time lapse (Figure 2C). The number  $n$  of cells that is appropriate to use might be different depending on the specifics of the developing tissue and its cellular architecture. In the case of the zebrafish PSM we found little difference in the resulting AGETs depending on whether these were calculated using 10, 5 or 1 cell-neighbourhoods (Figures 2B and 2C). The averaging method used (mean or median) had negligible effect in this case (Figures 2B and 2C). Altogether this suggests that AGETs calculated for cells in the developing zebrafish PSM are very robust to the methods used to calculate them (Figure 2). The resulting tissue-level patterns were also found to be robust to the number of cells and the averaging method used (Figure S3). Once AGETs have been calculated for every cell in the time lapse, we can visualise the resulting pattern at the level of the tissue by plotting the positions of the cells over time and colour-coding them according to their approximated gene expression values (Figures 2D and Videos S1 and S2). We refer to these datasets as *in silico reporters* since they allow us to visualise gene expression dynamics over the course of development (Videos S1 and S2). The resulting tissue-level pattern can be quantified and compared to that of the HCRs, revealing a good match between the two (Figure 2D).

**Approximated gene expression trajectories can be used 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 *tbxta*, *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. 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 minimizes the difference between the target and the simulated data. We chose to use a Markov Chain Monte Carlo (MCMC) algorithm<sup>31</sup> to use as our parameter sampling method since MCMC has been extensively used and repeatedly validated for GRN inference.<sup>32</sup> In addition, MCMC and has the advantage of providing a population of candidate networks by approximating the entire posterior distribution for each GRN parameter.

**Optimal fits are obtained using 10% of all approximated gene expression trajectories for fitting**

We first sought to identify the optimal number of AGETs to fit to: one that is computationally feasible while producing high quality fits. We fit to various numbers of AGETs, between 10 and 200, and computed the overall quality of the resulting fits for generated parameter sets. We found that while good fits can be obtained when fitting to as few as 10 AGETs, fitting to increased numbers of AGETs increases the proportion of inferred well-fitting parameter sets.



**Figure 3. Goodness of fits and performance of the GRN corresponding to the maximum a posteriori (MAP) parameters when fitting to an optimum of 200 AGETs**

(A) Likelihood score improves (increases) with the number of AGETs used for model fitting. GRNs were fit using 10, 50, 100, and 200 AGETs (of a total of around 1903). To compute overall model fit when fitting to different numbers of AGETs, we computed the likelihood using all 1903 AGETs for the final 2000 parameter sets of each run (see Materials and Methods). Fitting to 200 cells produced the likelihood distributions with the best mean likelihood values overall, corresponding to the best fitting parameter sets. All GRNs presented from here on were obtained by fitting to 200 AGETs.

(B) Cell tracks whose AGETs have been used for model fitting (200 cell tracks shown in thin red lines), used for model fitting where the fits are shown in D. (3 cell tracks shown in thick red lines) and used for model testing but not fitting, and where the fits are shown in E. (9 cell tracks shown in thick black lines). Cell tracks are shown against the outline of the PSM to convey their relative position within the tissue.

**Figure 3. Continued**

(C) Topology of the best fitting network: this is the parameter set with the lowest likelihood score when computed on all simulated AGETs. Positive interactions are blue arrows, negative interactions are red T-bars, and small/zero interaction values are indicated by gray lines. The magnitude of the interaction is indicated by color intensity. Parameter values are shown in [Table S1](#) (Materials and Methods) (D) Comparison of simulated vs. target AGETs for three cells used in the fitting procedure, which are initially located at i. 12%, ii. 23% and iii. 27% PA within the PSM. Red: *tbxta* expression, yellow: *tbx16* expression, and blue: *tbx6* expression. Thick lines denote simulated gene expression and thin lines, the target AGETs.

(E) Comparison of simulated vs. target AGETs for nine cells not used in the fitting procedure, which are initially located at i. 5%, ii. 11%, iii. 12%, iv. 19%, v. 28%, vi. 32%, vii. 36%, viii. 36%, and ix. 47% PA within the PSM. Red: *tbxta* expression, yellow: *tbx16* expression, and blue: *tbx6* expression. Thick lines denote simulated gene expression and thin lines, the target AGETs.

(F) Snapshots showing simulated *tbx* expression quantified along the PA axis of the PSM at three time frames in the time-lapse: i. 20, ii. 40 and iii. 60. Each dot represents the simulated value of a gene (*tbxta* in red, *tbx16* in yellow, and *tbx6* in blue) in a cell at a given time point. Lines represent maximum projection measures of *tbx* patterns from the HCRs (smoothed average gene expression profiles). Note that the anterior down-regulation of *tbx6* is not recapitulated; this was intentionally not fit to as factors not included in the current GRN formulation are responsible for this feature of the pattern.

Model fits will improve with the number of AGETs used as reflected by the increase in the mean likelihood ([Figure 3A](#)) and in the mean acceptance proportion for each run ([Figure S5A](#)), decrease in likelihood spread ([Figure 3A](#)) and leveling in the mean auto-correlation score for each run as observed when using 200 AGETs for fitting or more ([Figure S5B](#)). Since the goodness of the fits and the convergence of the MCMC algorithm will increase up to 200 AGETs, plateauing thereafter ([Figure 3A](#) and [S5](#)) we settle on 200 AGETs (approximately 10% of all available AGETs) as the optimal number of AGETs with which to reverse-engineer GRNs for the rest of this study.

The 200 AGETs used for model fitting were selected randomly and are distributed uniformly throughout the PSM ([Figure 3B](#), thin red lines). We only selected AGETs from cells that had been consecutively tracked for the entire duration of the time lapse (61 frames). We expect that the optimal number of AGETs required to obtain good fits will be system-specific.

**The maximal posterior probability network recapitulates *tbx* gene expression dynamics at the cellular and tissue levels**

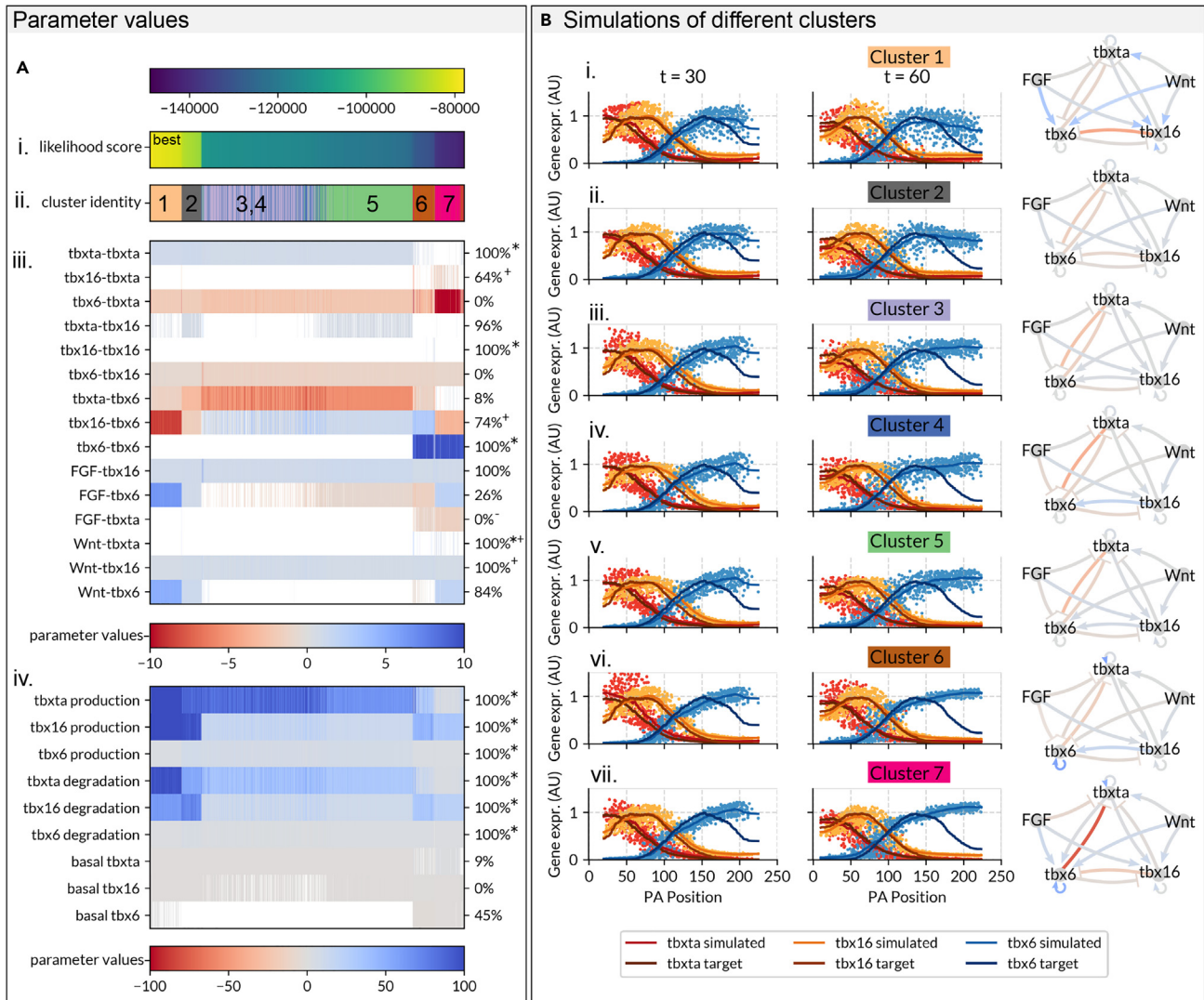
MCMC inference yields a collection of parameter sets (samples) that together approximate the posterior distribution of the GRN's parameters. For every parameter, we obtain a probability distribution across its possible values, which provides information about the values that are most likely to produce good fits. We first chose to explore the network corresponding to the parameter set with the overall highest posterior probability score: the network corresponding to the maximum a posteriori - or MAP - sample ([Figure 3C](#)). We use this network to simulate all 1903 available AGETs, and visualize the simulation on the tracks ([Videos S2](#) and [S3](#)). We validate the quality of the inferred network by both comparing single AGETs with their simulated counterparts ([Figures 3D](#) and [3E](#)), and by comparing the whole tissue-level gene expression profiles over time ([Figure 3F](#)). When simulating single AGETs we find that the MAP network recapitulates both AGETs used for fitting and not used for fitting alike ([Figures 3B–3D](#)). 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 recover AGET gene expression levels and trends.

We are especially interested in how well the simulations recapitulate whole tissue patterning dynamics, as these emerge mostly from simulating AGETs that have not been used for model training (approximately 90% of all AGETs). [Figure 3F](#) i-iii. shows simulated T-*bx* expression for each cell along the normalized posterior to anterior axis of the PSM (dots) at time frames 20, 40, and 60. Simulated data have been fit at each separate time point by curves which are then normalized (darker curves) and compared to the curves previously obtained from the AGETs (shown as lighter curves) ([Figure 3F](#)). Overall, simulations recapitulate target tissue-level gene expression very well (snapshots in [Figure 3F](#) and full simulations in [Videos S2](#) and [S3](#)).

Note that there is a discrepancy between the AGETs and the simulated anterior *tbx6* expression. The formulated GRN is unrealistic in this region, as additional factors secreted from the somites are known to down-regulate this transcription factor.<sup>33</sup> For this reason we intentionally excluded this region during the parameter optimization procedure by omitting cells in the anterior-most region of the PSM, and accordingly, simulated gene expression here is inaccurate. 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.<sup>26</sup> 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 biology and might suggest that the three genes considered are indeed causally responsible for much of this patterning system.

**Reverse engineering gene regulatory networks using 200 approximated gene expression trajectories yields seven clusters of solutions**

MCMC is a parameter sampling algorithm, and as such it will return an approximated posterior distribution for the GRN parameters instead of a single point estimate. This provides a range of candidate networks that can be subsequently analyzed 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 are not, taking instead a broad range of values across the inferred networks. Such information can lead to interesting hypotheses regarding which aspects of the pattern evolution might be most strongly acting on.



**Figure 4. Seven main clusters of GRN topologies that recapitulate Tbx expression along the developing zebrafish PSM**

Reverse-engineering GRN topology using sets of 200 AGETs yields seven main clusters of solutions which recapitulate *tbx* spatiotemporal gene expression along the PSM, as well as known genetic interactions.

(A) The likelihood score for the final converging 2000 sets of parameters is computed on all the tracks. Networks are displayed in order of increasing likelihood (improved fit). A.i. Likelihood scores of successful parameter sets. A.ii. Cluster identity of parameter sets, from k-means clustering. A.iii. Parameters representing interactions between GRN nodes. A.iv. Parameters representing production, degradation and basal gene expression. Parameters indicated by an asterisk (\*) are set as positive priors when fitting. E. Resulting simulations from representative networks from each cluster.

While in the previous section we analyzed 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 (Figure 4). To do this, we calculated the likelihood score across all the tracks (Materials and Methods) for the final 2000 parameter sets for every fitting run which reached convergence. As the MCMC algorithm used can produce small numbers of poorly fitting parameter sets, even when convergence has been reached, we then excluded all parameter sets with an overall likelihood score of less than  $-15000$ , as well as parameter sets with unrealistic values (genetic interactions  $< 100$  or  $> -100$ ). This left us with a total of 3600 candidate networks. We plotted these networks in order of increasing likelihood (better fit) (Figures 4A–4D) color coding according to whether a given interaction is positive (activation, red) or negative (repression, blue) to visualize the different predicted network topologies.

We clustered these networks using k-means clustering on scaled parameter values (Materials and Methods). We initially clustered them into 13 clusters in total (as determined by an elbow plot - supplemental data), and considered only those clusters which included 5% or more of the total number of networks, reducing the number of clusters to seven. All clusters produce well-fitting patterns (Figure 4B).

Furthermore, since all seven clusters of networks have very similar likelihood scores, we treat them all as equally probable candidates until further biological insight helps us to discriminate between them.<sup>26</sup> We are confident that our methodology is effective at producing a range of topologically distinct and well fitting parameter sets, as evidenced by the accuracy of the resulting fits for all of the networks.

We investigated whether these networks recovered genetic interactions known from the literature. For instance, Wnt is known to directly activate *tbx16* in the zebrafish tailbud: Bouldin et al. identify a Wnt-binding promoter that drives *tbx16* expression.<sup>34</sup> All our parameter sets predict this interaction is positive, thus being supported by the literature (Figure 4A iii.). *tbx16* is predicted to activate *tbx6* in 74% of the inferred networks (Figure 4A iii.). The nature of this interaction has been experimentally demonstrated using heat-shock transgenic lines to up-regulate *tbx16* expression which led to an increased expression of *tbx6*.<sup>34</sup> *tbx16* mutants have a loss of *tbx6* expression,<sup>35</sup> further corroborating this interaction. This work also<sup>34</sup> also shows that over-expression of *tbx16* results in a decrease of *tbxta* in the tailbud, but that this phenotype is rescued by over-expression of Wnt, suggesting that this is an indirect interaction occurring via Wnt. 64% of our parameter sets have a positive interaction between *tbx16* and *tbx6* (Figure 4A iii.). Similarly,<sup>36</sup> show that at the tailbud stage, loss of FGF causes an expansion of *tbxta*, suggesting that FGF inhibits *tbxta*; in all of our parameter sets, FGF inhibits *tbxta*. It is also known that FGF activates *tbx16* expression as dominant negative FGFR1 embryos display a loss of *tbx16* expression<sup>36, ?</sup>; again, in all of our parameter sets FGF is activating *tbx16* (Figure 4A iii.). Together, these data reveal a high degree of consensus between our parameter values and known interactions established in the literature.

### Reverse-engineered gene regulatory networks qualitatively recapitulate known results of experimental perturbations

In the previous section we validated our parameter sets and GRN inference methodology by comparing our parameter sets to known genetic interactions from the literature. To further validate our parameter sets, we attempted to replicate the effects of changing FGF and Wnt signaling. Bouldin et al.<sup>34</sup> over-express Wnt in the zebrafish tailbud by over-expressing beta-catenin under a heat-shock promoter. Four hours post heat-shock, the *tbx16* expression domain is expanded toward the anterior PSM as a result of an up-regulation of *tbx16* in this region. We replicated this experiment in silico by setting Wnt expression to 1.5 and maintained it for the duration of the simulation. We then ran these simulations on the tracks as before. When this simulation was run using the MAP network we find that at the final time-point of simulation *tbx16* is up-regulated in the anterior PSM, as in the experimental data (Figures 5A and 5B). We see, however, a down-regulation of *tbxta* which is opposite to the up-regulation of *tbxta* reported by Bouldin et al. This could be related to the fact that our model formulation does not include the known feedback loop between Wnt and *tbxta*.<sup>37</sup> When we run this in silico experiment with all the networks we find that the vast majority predict an average increase in *tbx16* expression (Figure 5C).

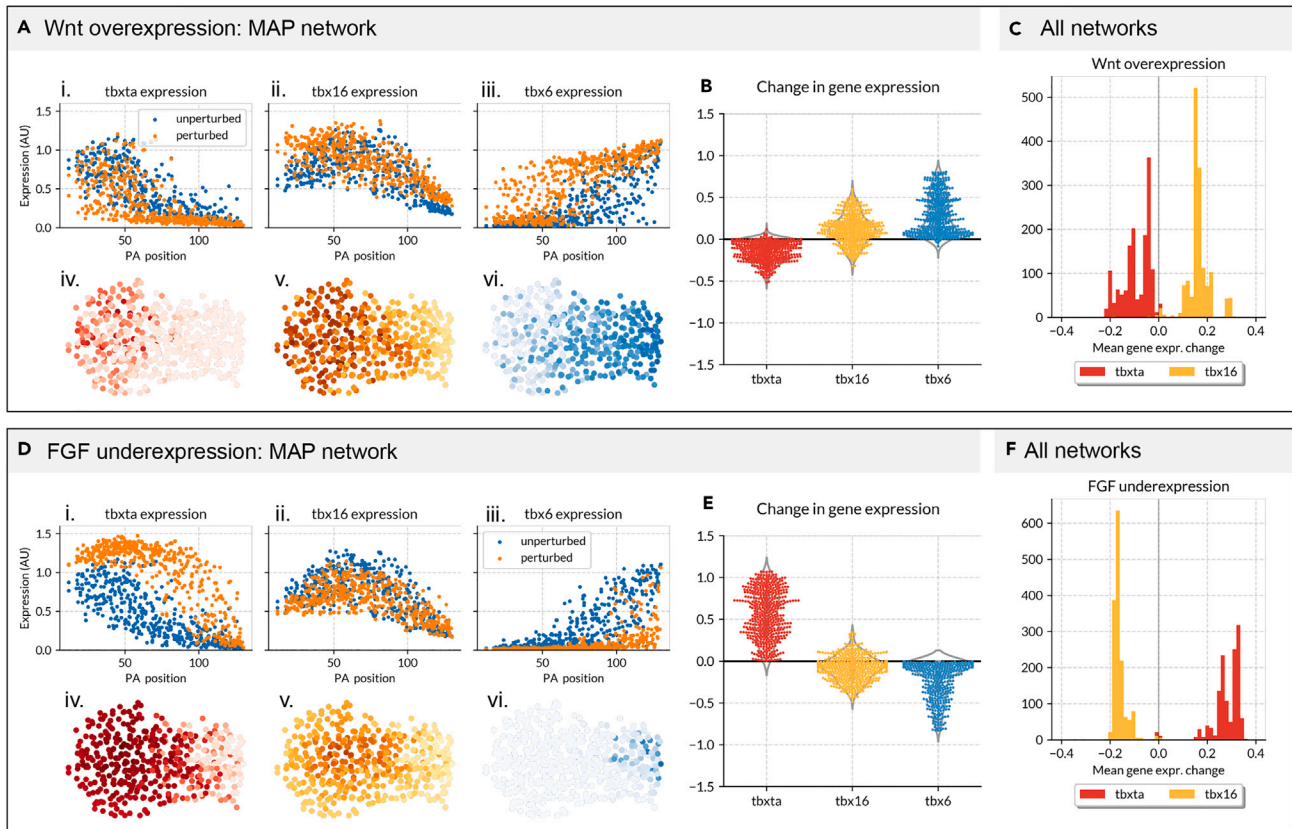
Goto et al. downregulate FGF expression by inducing a heat-shock dominant negative form of FGF at 8 ss.<sup>36</sup> At this time point, loss of FGF results in an expansion of *tbxta* expression 3 h after heatshock, at approximately the 15ss. We reproduced this experiment by setting FGF to 0.01 in every cell and simulating as before. For the MAP network, *tbxta* is significantly up-regulated after the downregulation of FGF, in agreement with the experimental results. *tbx16* is slightly down-regulated, and *tbx6* is significantly down-regulated (Figures 5D and 5E). When we repeated this analysis on all successful parameter sets, a majority recapitulate an increase in *tbxta* expression, as represented by the average difference between *tbxta* expression in each cell in WT and perturbed simulations (Figure 5F).

### 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 toward 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.

Here, we propose a method to reverse-engineer gene regulatory networks using datasets that are relatively straightforward to generate in an increasing number of model and non-model species spanning the range of animal phylogeny. This will make it possible to construct AGETs and therefore 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 general principles that were inaccessible with previous approaches. Furthermore, this methodology will find applications well-beyond beyond the study of developmental evolution. In particular, we anticipate technique will be particularly useful in 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 organization is paramount and has proven not at all trivial.

Finally, our methodology for the construction of AGETs provides a way in which to visualize approximated gene expression dynamics and patterns in the form of in-silico reporters, not with the aim of replacing live reporters, but rather by providing an approximation until live reporters become available. There is in principle no limit to the number of genes that can be reported by an in silico reporter line which could be used for hypothesis generation and to compare the relative co-expression of previously unexplored combinations of genes. In silico reporter lines can also be readily extended to non-model organisms where transgenic reporters are not yet available. All in all, we expect that this methodology will find a broad range of applications in developmental evolution and beyond, contributing to advance the data-driven study of patterning dynamics.



**Figure 5. GRNs reverse-engineered using 200 AGETs recapitulate known results of experimental perturbations**

(A) Results of simulating Wnt over-expression using the MAP network. (A.i) *tbxta*, (A.ii) *tbx16*, and (A.iii) *tbx6* expression in unperturbed (blue) and perturbed (orange) simulations using the MAP network. X axis denotes the posterior to anterior position along the PSM. Y axis represents gene expression levels in arbitrary units. Each dot represents the expression level in a cell at the last time point of the simulation. A.iv. *tbxta* (red), A.v. *tbx16* (yellow), and (A.vi) *tbx6* (blue) expression shown in cells within the PSM. Posterior left, dorsal up.

(B) Difference in gene expression between the unperturbed and perturbed simulations calculated for every cell at the last time point of the simulations.

(C) The mean gene expression difference in single cells between perturbed and unperturbed simulations, calculated for every network and displayed as a histogram for *tbxta* (red) and *tbx16* (yellow).

(F) Results of simulating FGF under-expression.

(D.i). *tbxta*, (D.ii) *tbx16* and (D.iii) *tbx6* expression in unperturbed (blue) and perturbed (orange) simulations using the MAP network. X axis denotes the posterior to anterior position along the PSM. Y axis represents gene expression levels in arbitrary units. Each dot represents the expression level in a cell at the last time point of the simulation. (D.iv) *tbxta* (red), (D.v) *tbx16* (yellow), and (D.vi) *tbx6* (blue) expression is shown in cells within the PSM. Posterior left, dorsal up.

(E) Difference in gene expression between the unperturbed and perturbed simulations calculated for every cell at the last time point of the simulations.

(F) The mean gene expression difference in single cells between perturbed and unperturbed simulations, calculated for every network and displayed as a histogram for *tbxta* (red) and *tbx16* (yellow).

### Limitations of the study

While our methodology is in principle applicable across a wide variety of developmental processes, it would have to be substantially modified to accommodate highly dynamic cases of pattern formation. The AGET construction algorithm makes the assumption that a cell's gene expression is strongly correlated with its position within a tissue. This is an assumption that applies well to Tbox patterning in the zebrafish PSM, where cells are known to express the gene sequence *tbxta*, *tbx16*, *tbx6* as they differentiate and transit toward the somites.<sup>26</sup> However, this assumption does not always apply; for example in the case of oscillatory gene expression in during somitogenesis<sup>38</sup> or neural differentiation<sup>39</sup> which are highly dynamic processes where gene expression in a cell might change before the cell's position. In such systems, our current AGET construction methods would need to be revised and HCRs would need to cover the dynamics of the process at a much higher resolution. Still, since this method is based on approximating dynamics from static images, one would risk missing very rapid or transient gene expression dynamics and in such cases we would recommend that the resulting AGETs be validated using transgenic reporter lines that are known to faithfully recapitulate the gene expression dynamics. This was unfortunately not possible for the Tbox genes as existing reporter lines are tagged with long-lived EGFP, making the reported expression dynamics unsuitable for AGET validation.<sup>40</sup> To mitigate this in the case of the Tbox genes, we have performed extensive model

validation using features of the GRN that have not been fit to, such as GRN topology and the ability to predict perturbation experiments.

## RESOURCE AVAILABILITY

### Lead contact

Further information and requests for resources and reagents should be directed to and will be fulfilled by the lead contact, Berta Verd ([berta.verdfernandez@biology.ox.ac.uk](mailto:berta.verdfernandez@biology.ox.ac.uk)).

### Materials availability

This study did not generate new unique reagents.

### Data and code availability

- Data and code available from <https://github.com/Shannon-E-Taylor/AGETS/tree/main>.
- All original code has been deposited at Zenodo and is publicly available as of the date of publication. DOIs are listed in the [key resources table](#).
- Microscopy data reported in this article was originally reported and described in.<sup>27</sup>
- Any additional information required to reanalyze the data reported in this article is available from the [lead contact](#) on request.

## ACKNOWLEDGMENTS

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

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 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. Y.W. is supported by a summer vacation stipend from St Catharine's College, University of Cambridge. B.P. was supported by the Alan Turing Institute and University College London. 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). B.V. was funded by ERC StG COUNTS 101163722.

## AUTHOR CONTRIBUTIONS

Conceptualization: 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.

## DECLARATION OF INTERESTS

The authors declare no competing interests.

## STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- [KEY RESOURCES TABLE](#)
- [EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS](#)
  - Animal lines and husbandry
- [METHOD DETAILS](#)
  - *In Situ* Hybridisation Chain Reaction (HCR)
  - Immunohistochemistry
  - Imaging and image analysis
  - Image analysis processing pipeline for HCR images
  - Aligning point clouds with ICP
  - AGET construction
  - Mathematical model formulation
  - Model fitting: MCMC approach
  - Calculating a likelihood score
  - Filtering and clustering

## SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2024.110840>.

Received: February 6, 2024

Revised: March 20, 2024

Accepted: August 21, 2024

Published: August 30, 2024

**REFERENCES**

- Reinitz, J., and Sharp, D.H. (1996). Gene circuits and their uses. In *Integrative Approaches to Molecular Biology* (MIT Press), pp. 253–272.
- Liang, S., Fuhrman, S., and Somogyi, R. (1998). Reveal, a general reverse engineering algorithm for inference of genetic network architectures. In *Pac. Symp. Biocomput., 3Pac. Symp. Biocomput.* (Citeseer), pp. 18–29.
- D'haeseleer, P., Liang, S., and Somogyi, R. (2000). Genetic network inference: from co-expression clustering to reverse engineering. *Bioinformatics* 16, 707–726.
- Gardner, T.S., and Faith, J.J. (2005). Reverse-engineering transcription control networks. *Phys. Life Rev.* 2, 65–88.
- Rockman, M.V. (2008). Reverse engineering the genotype–phenotype map with natural genetic variation. *Nature* 456, 738–744.
- He, F., Balling, R., and Zeng, A.-P. (2009). Reverse engineering and verification of gene networks: principles, assumptions, and limitations of present methods and future perspectives. *J. Biotechnol.* 144, 190–203.
- Jaeger, J., and Monk, N.A. (2010). Reverse engineering of gene regulatory networks. *Learning and inference in computational systems biology* 9, 34.
- Crombach, A., Wotton, K.R., Cicin-Sain, D., Ashyraliyev, M., and Jaeger, J. (2012). Efficient reverse-engineering of a developmental gene regulatory network. *PLoS Comput. Biol.* 8, e1002589.
- Verd, B., Crombach, A., and Jaeger, J. (2017). Dynamic maternal gradients control timing and shift-rates for drosophila gap gene expression. *PLoS Comput. Biol.* 13, e1005285.
- Verd, B., Clark, E., Wotton, K.R., Janssens, H., Jiménez-Guri, E., Crombach, A., and Jaeger, J. (2018). A damped oscillator imposes temporal order on posterior gap gene expression in drosophila. *PLoS Biol.* 16, e2003174.
- Manu, S., Reinitz, J., Surkova, S., Spirov, A.V., Gursky, V.V., Janssens, H., Kim, A.R., Radulescu, O., Vanario-Alonso, C.E., Sharp, D.H., and Samsonova, M. (2009). Canalization of gene expression and domain shifts in the drosophila blastoderm by dynamical attractors. *PLoS Comput. Biol.* 5, e1000303.
- Balaskas, N., Ribeiro, A., Panovska, J., Dessaud, E., Sasai, N., Page, K.M., Briscoe, J., and Ribes, V. (2012). Gene regulatory logic for reading the sonic hedgehog signaling gradient in the vertebrate neural tube. *Cell* 148, 273–284.
- Verd, B., Monk, N.A., and Jaeger, J. (2019). Modularity, criticality, and evolvability of a developmental gene regulatory network. *Elife* 8, e42832.
- Jaeger, J., Blagov, M., Kosman, D., Kozlov, K.N., Sharp, D.H., Reinitz, J., Myasnikova, E., Myasnikova, E., Surkova, S., Vanario-Alonso, C.E., and Samsonova, M. (2004). Dynamical analysis of regulatory interactions in the gap gene system of drosophila melanogaster. *Genetics* 167, 1721–1737.
- Cohen, M., Kicheva, A., Ribeiro, A., Blassberg, R., Page, K.M., Barnes, C.P., and Briscoe, J. (2015). Ptch1 and gli regulate shh signalling dynamics via multiple mechanisms. *Nat. Commun.* 6, 6709–6712.
- Kicheva, A., Bollenbach, T., Ribeiro, A., Valle, H.P., Lovell-Badge, R., Episkopou, V., and Briscoe, J. (2014). Coordination of progenitor specification and growth in mouse and chick spinal cord. *Science* 345, 1254927.
- El-Sherif, E., Zhu, X., Fu, J., and Brown, S.J. (2014). Caudal regulates the spatiotemporal dynamics of pair-rule waves in tribolium. *PLoS Genet.* 10, e1004677.
- Wu, H., Jiao, R., and Ma, J. (2015). Temporal and spatial dynamics of scaling-specific features of a gene regulatory network in drosophila. *Nat. Commun.* 6, 1–13.
- Averbukh, I., Lai, S.-L., Doe, C.Q., and Barkai, N. (2018). A repressor-decay timer for robust temporal patterning in embryonic drosophila neuroblast lineages. *Elife* 7, e38631.
- Schröter, C., Ares, S., Morelli, L.G., Isakova, A., Hens, K., Soroldoni, D., Gajewski, M., Jülicher, F., Maerkl, S.J., Deplancke, B., and Oates, A.C. (2012). Topology and dynamics of the zebrafish segmentation clock core circuit. *PLoS Biol.* 10, e1001364.
- Rayon, T., Stamatakis, D., Perez-Carrasco, R., Garcia-Perez, L., Barrington, C., Melchionda, M., Exelby, K., Lazaro, J., Tybulewicz, V.L.J., Fisher, E.M.C., and Briscoe, J. (2020). Species-specific pace of development is associated with differences in protein stability. *Science* 369, eaba7667.
- Uriu, K., Morishita, Y., and Iwasa, Y. (2010). Random cell movement promotes synchronization of the segmentation clock. *Proc. Natl. Acad. Sci. USA* 107, 4979–4984.
- Crombach, A., Wotton, K.R., Jiménez-Guri, E., and Jaeger, J. (2016). Gap gene regulatory dynamics evolve along a genotype network. *Mol. Biol. Evol.* 33, 1293–1307.
- Kicheva, A., Cohen, M., and Briscoe, J. (2012). Developmental pattern formation: insights from physics and biology. *Science* 338, 210–212.
- Huch, M., Knoblich, J.A., Lutolf, M.P., and Martínez-Arias, A. (2017). The hope and the hype of organoid research. *Development* 144, 938–941.
- Fulton, T., Speiss, K., Thomson, L., Wang, Y., Clark, B., Hwang, S., Paige, B., Verd, B., and Steventon, B. (2022). Cell rearrangement generates pattern emergence as a function of temporal morphogen exposure. Preprint at bioRxiv. <https://doi.org/10.1101/2021.02.05.429898>.
- Thomson, L., Muresan, L., and Steventon, B. (2021). The zebrafish presomitic mesoderm elongates through compaction-extension. *Cells & Development* 168, 203748.
- Mongera, A., Rowghanian, P., Gustafson, H.J., Shelton, E., Kealhofer, D.A., Carn, E.K., Serwane, F., Lucio, A.A., Giammona, J., and Campàs, O. (2018). A fluid-to-solid jamming transition underlies vertebrate body axis elongation. *Nature* 561, 401–405.
- Choi, H.M.T., Schwarzkopf, M., Fornace, M.E., Acharya, A., Artavanis, G., Stegmaier, J., Cunha, A., and Pierce, N.A. (2018). Third-generation in situ hybridization chain reaction: multiplexed, quantitative, sensitive, versatile, robust. *Development* 145, dev165753.
- Rusinkiewicz, S., and Levoy, M. (2001). Efficient variants of the icp algorithm. In *Proceedings third international conference on 3-D digital imaging and modeling (IEEE)*, pp. 145–152.
- Foreman-Mackey, D., Hogg, D.W., Lang, D., and Goodman, J. (2013). Emcee: The Mcmc Hammer, 125 (Publications of the Astronomical Society of the Pacific), p. 306.
- Ram, R., and Chetty, M. (2009). Mcmc based bayesian inference for modeling gene networks. In *IAPR International Conference on Pattern Recognition in Bioinformatics (Springer)*, pp. 293–306.
- Kawamura, A., Koshida, S., Hijikata, H., Ohbayashi, A., Kondoh, H., and Takada, S. (2005). Groucho-associated transcriptional repressor ripply1 is required for proper transition from the presomitic mesoderm to somites. *Dev. Cell* 9, 735–744.
- Bouldin, C.M., Manning, A.J., Peng, Y.-H., Farr, G.H., Hung, K.L., Dong, A., and Kimelman, D. (2015). Wnt signaling and tbx16 form a bistable switch to commit bipotential progenitors to mesoderm. *Development* 142, 2499–2507.
- Fior, R., Maxwell, A.A., Ma, T.P., Vezzano, A., Moens, C.B., Amacher, S.L., Lewis, J., Saúde, L., and Saude, L. (2012). The differentiation and movement of presomitic mesoderm progenitor cells are controlled by mesogenin 1. *Development (Cambridge, England)* 139, 4656–4665. <https://doi.org/10.1242/dev.078923>.
- Goto, H., Kimmey, S.C., Row, R.H., Matus, D.Q., and Martin, B.L. (2017). Fgf and canonical wnt signaling cooperate to induce paraxial mesoderm from tailbud neuromesodermal progenitors through regulation of a two-step epithelial to mesenchymal transition. *Development* 144, 1412–1424.
- Martin, B.L., and Kimelman, D. (2008). Regulation of canonical wnt signaling by brachyury is essential for posterior mesoderm formation. *Dev. Cell* 15, 121–133.
- Soroldoni, D., Jörg, D.J., Morelli, L.G., Richmond, D.L., Schindelin, J., Jülicher, F., and Oates, A.C. (2014). A doppler effect in embryonic pattern formation. *Science* 345, 222–225.
- Kaise, T., and Kageyama, R. (2021). Hes1 oscillation frequency correlates with activation of neural stem cells. *Gene Expr. Patterns* 40, 119170.
- Ban, H., Yokota, D., Otosaka, S., Kikuchi, M., Kinoshita, H., Fujino, Y., Yabe, T., Ovara, H., Izuka, A., Akama, K., et al. (2019). Transcriptional autoregulation of zebrafish tbx6 is required for somite segmentation. *Development* 146, dev177063.
- Moro, E., Ozhan-Kizil, G., Mongera, A., Beis, D., Wierzbicki, C., Young, R.M., Bournele, D., Domenichini, A., Valdivia, L.E., Lum, L., et al. (2012). In vivo wnt signaling tracing through a transgenic

- biosensor fish reveals novel activity domains. *Dev. Biol.* 366, 327–340.
42. Zhou, Q.-Y., Park, J., and Koltun, V. (2018). Open3D: A modern library for 3D data processing. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1801.09847>.
  43. Kimmel, C.B., Ballard, W.W., Kimmel, S.R., Ullmann, B., and Schilling, T.F. (1995). Stages of embryonic development of the zebrafish. *Dev. Dynam.* 203, 253–310.
  44. Hirsinger, E., and Steventon, B. (2017). A versatile mounting method for long term imaging of zebrafish development. *JoVE* 119, e55210.
  45. Mjolsness, E., Sharp, D.H., and Reinitz, J. (1991). A connectionist model of development. *J. Theor. Biol.* 152, 429–453.
  46. Collier, J.R., Monk, N.A., Maini, P.K., and Lewis, J.H. (1996). Pattern formation by lateral inhibition with feedback: a mathematical model of delta-notch intercellular signalling. *J. Theor. Biol.* 183, 429–446.
  47. Verd, B., Crombach, A., and Jaeger, J. (2014). Classification of transient behaviours in a time-dependent toggle switch model. *BMC Syst. Biol.* 8, 43.

## STAR★METHODS

### KEY RESOURCES TABLE

| REAGENT or RESOURCE                                   | SOURCE   | IDENTIFIER  |
|---|--|---|
| <b>Antibodies</b>                                     |  |   |
| Diphosphorylated Mouse ERK antibody                   | Sigma  | M9692-200UL   |
| Secondary Mouse Alexa 647nm antibody                  | Invitrogen   | A21236  |
| <b>Chemicals, peptides, and recombinant proteins</b>  |  |   |
| 4% paraformaldehyde (PFA)                             | Sigma  | CAS no: 30,525-89-4   |
| Agarose, low gelling temperature                      | Sigma  | A9414   |
| Dulbecco's Phosphate Buffered Saline (PBS)            | Sigma  | D8537-500ML   |
| DAPI  | Sigma  | CAS no: 28718-90-3  |
| EDTA  | Sigma  | CAS no: 60-00-4   |
| Triton X-   | Sigma  | CAS no:9002-93  |
| Fetal bovine serum (heat-inactivated)                 | ThermoFisher Scientific                            | 10437028  |
| Bovine serum Albumin                                  | Sigma  | A7906-10G   |
| VECTASHIELD Antifade mounting medium                  | Vector Laboratories                                | H-1000-10   |
| Tween 20  | ThermoFisher Scientific                            | AAJ20605A   |
| Methylcellulose                                       |  | M0512   |
| SSC buffer  | Scientific Laboratory Supplies                     | S6639-1L<br>S6639-1L  |
| <b>Critical commercial assays</b>                     |  |   |
| <i>In situ</i> HCR v3.0                               | Molecular Instruments                              | N/A   |
| <b>Deposited data</b>                                 |  |   |
| TBOX and FGF/Wnt gene expression in nuclei            | This paper   | <a href="https://github.com/Shannon-E-Taylor/AGETS/tree/main/01_AGET_construction/Source_images">https://github.com/Shannon-E-Taylor/AGETS/tree/main/01_AGET_construction/Source_images</a>                                     |
| Cell tracking data                                    | This paper   | <a href="https://github.com/Shannon-E-Taylor/AGETS/blob/main/01_AGET_construction/Tracking_data/Tracks_M4_Kay.csv">https://github.com/Shannon-E-Taylor/AGETS/blob/main/01_AGET_construction/Tracking_data/Tracks_M4_Kay.csv</a> |
| AGETs code  | This paper   | <a href="https://doi.org/10.5281/zenodo.12607614">https://doi.org/10.5281/zenodo.12607614</a>   |
| <b>Experimental models: Organisms/strains</b>         |  |   |
| Wildtype Zebrafish embryos - Tüpfel long fin (TL)     | European Zebrafish Resource Center                 | ZDB-GENO-990623-2   |
| Wildtype Zebrafish embryos - AB                       | European Zebrafish Resource Center                 | ZDB-GENO-960809-7   |
| The Tg(7xTCF- Xla.Sia:GFP)                            | Steven Wilson Lab; Moro et al., 2012 <sup>41</sup> | ZDB-TGCONSTRCT-110113-1   |
| <b>Software and algorithms</b>                        |  |   |
| Imaris (v9.2.1)                                       | Bitplane   | <a href="https://imaris.oxinst.com/">https://imaris.oxinst.com/</a>   |
| Python v3.7   | Python Software Foundation                         | <a href="https://www.python.org/">https://www.python.org/</a>   |
| Open3D (v0.13.0)                                      | Zhou et al. 2018 <sup>42</sup>                     | <a href="http://www.open3d.org/">http://www.open3d.org/</a>   |
| Emcee (v3.1.0)  | Foreman-Mackey et al. 2012 <sup>31</sup>           | <a href="https://emcee.readthedocs.io/en/stable/">https://emcee.readthedocs.io/en/stable/</a>   |
| Anaconda environment detailing other package versions | This paper   | <a href="https://github.com/Shannon-E-Taylor/AGETS/blob/main/environment.yml">https://github.com/Shannon-E-Taylor/AGETS/blob/main/environment.yml</a>   |
| sklearn v0.24.2                                       | Scikit-learn                                       | <a href="https://scikit-learn.org/stable/">https://scikit-learn.org/stable/</a>   |

(Continued on next page)

**Continued**

| REAGENT or RESOURCE            | SOURCE | IDENTIFIER  |
|--------------------------------|--------|---|
| Other                          |        |   |
| Inverted Confocal microscope   | Zeiss  | Zeiss LSM700  |
| Two-photon confocal microscope |        | TriM Scope II Upright<br>2-photon scanning fluorescence microscope<br>equipped Insight DeepSee dual-line laser<br>(tunable<br>710-1300 nm fixed 1040 nm line) |
| 35 mm glass bottom dish        | MatTek | P35G-1.5-10-C   |

**EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS****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 Tüpfel Long Fin (TL), AB or AB/TL. The Tg(7xTCF-Xla.Sia:GFP) reporter line<sup>41</sup> was provided by the Steven Wilson laboratory. Embryos were staged as in.<sup>43</sup>

**METHOD DETAILS*****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.<sup>29</sup> 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 1 h at room temperature. Our readout 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<sup>27</sup>). The developing embryo was imaged with a 25×1.05 NA water dipping objective. Embryos were positioned laterally in low melting agarose with the entire tail cut free to allow for normal development.<sup>44</sup> Tracks were generated automatically and validated manually using the Imaris imaging software.

**Image analysis processing pipeline for HCR images**

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 (Figure S1B). 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 (Figure S1C). Background noise in the data was reduced by setting lower-bound thresholds for every gene. These thresholds were chosen such that *tbxta* and *tbx16* would appear restricted to the posterior end of the PSM (Figure S1Di, 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 (Figure S1Diii 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 (Figure 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 *tbxta*, *tbx16* and *tbx6* levels were extracted (Figure S1E). The same pipeline was used to obtain the levels of signals Wnt and FGF in single cells.

### Aligning point clouds with ICP

We used the Python library Open3d<sup>42</sup> and the implementation of the point-to-plane ICP (Iterative Closest Point) algorithm therein<sup>30</sup> 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.<sup>8,14,45</sup>

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

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

where  $R_a$  and  $\lambda_a$  respectively represent the rates of mRNA production and decay.  $\phi$  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:

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

where

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

$G = \{tbxta, tbx16, tbx6\}$  refers to the set of T-box genes while  $S = \{Wnt, FGF\}$  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.<sup>46,47</sup>

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<sup>31</sup> 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, we used a uniform prior from  $-200$  to  $+200$ , except when the priors were further restricted to positive values only, and fitted to the timescale used in the simulation. The timescale 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 96 walkers and for a total of 10'000 steps. We ran the MCMC independently at least 5 times for each parameter set, as this sometimes obtained different parameter distributions. Each fitting run was then extended for a further 10'000 iterations if it had not reached convergence. We defined convergence using the K-S test to check that the parameter distributions from the final 5000 parameter sets were statistically distinguishable from the  $-10000$ th to  $-8000$ th parameter sets, for at least 20 of the parameters. We used a significance threshold of 0.001.

### Calculating a likelihood score

To compute the log likelihood score during MCMC, we computed the squared distance between the target (AGET) and simulated gene expression for all cells fit to. This value was normalized using a custom  $\sigma$  value for each gene. More specifically, we used the following likelihood formulation:

$$L = -0.5 \cdot \sum_{a=0}^{a_{\max}} \sum_{c=0}^{c_{\max}} \sum_{t=0}^{t_{\max}} \frac{(T_{a t}^c - S_{a t}^c)^2}{\sigma_a}$$

Here,  $a$  refers to the gene being fit to,  $c$  refers to the cell being fit to; in this study  $c_{\max}$  was between 10 and 200.  $t$  refers to the time point of the simulation; in this study  $t_{\max} = 61$ .  $T$  refers to the Target - AGET - gene expression, and  $S$  refers to simulated gene expression.  $\sigma$  refers to the standard deviation which was unique for each gene.  $\sigma_{tbxta} = 0.2$ ,  $\sigma_{tbx16} = 0.2$ ,  $\sigma_{tbx6} = 0.1$ .

To compute overall model fit when fitting to different numbers of AGETs, we re-computed the likelihood using the same likelihood formulation but this time using all 1903 AGETs for the final 2000 parameter sets of each run. This allowed us to compare the goodness of the overall fit regardless of the number of AGETs used for fitting.

### Filtering and clustering

We filtered parameter sets for further analysis as follows. We firstly only analyzed parameter sets from runs that had reached convergence, and were fit to 200 AGETs. We then excluded unrealistically large parameter values, keeping values between  $\pm 100$ . We finally excluded parameter sets with a likelihood score less than  $-15000$ . Unfiltered parameter sets are presented in [Figure S6](#) and filtered in [Figure S7](#).

We performed k-means clustering of parameter values filtered as above, using the 'Kmeans' function of the 'sklearn' library v0.24.2. We chose the number of clusters to use using an Elbow Plot (FIG). We then removed extremely small clusters (less than 5% of all parameter sets) to obtain the seven clusters of parameter values presented in this paper.