

# Single-cell morphodynamical trajectories enable prediction of gene expression accompanying cell state change

**Authors:** Jeremy Copperman<sup>3\*</sup>, Ian C. Mclean<sup>1</sup>, Sean M. Gross, Young Hwan Chang<sup>1,2</sup>, Daniel M. Zuckerman<sup>1,2\*</sup>, and Laura M. Heiser<sup>1,2\*</sup>

## Affiliations:

<sup>1</sup> Department of Biomedical Engineering, Oregon Health and Science University, Portland OR 97239, U.S.A.

<sup>2</sup> Knight Cancer Institute, Oregon Health and Science University, Portland OR 97239, U.S.A

<sup>3</sup> Cancer Early Detection Advanced Research Center, Oregon Health and Science University, Portland OR 97239, U.S.A.

\* corresponding authors ([copperma@ohsu.edu](mailto:copperma@ohsu.edu), [zuckermd@ohsu.edu](mailto:zuckermd@ohsu.edu), [heiserl@ohsu.edu](mailto:heiserl@ohsu.edu))

## Abstract

Extracellular signals induce changes to molecular programs that modulate multiple cellular phenotypes, including proliferation, motility, and differentiation status. The connection between dynamically adapting phenotypic states and the molecular programs that define them is not well understood. Here we develop data-driven models of single-cell phenotypic responses to extracellular stimuli by linking gene transcription levels to “morphodynamics” – changes in cell morphology and motility observable in time-lapse image data. We adopt a dynamics-first view of cell state by grouping single-cell trajectories into states with shared morphodynamic responses. The single-cell trajectories enable development of a first-of-its-kind computational approach to map live-cell dynamics to snapshot gene transcript levels, which we term MMIST, Molecular and Morphodynamics-Integrated Single-cell Trajectories. The key conceptual advance of MMIST is that cell behavior can be quantified based on dynamically defined states and that extracellular signals alter the overall distribution of cell states by altering rates of switching between states. We find a cell state landscape that is bound by epithelial and mesenchymal endpoints, with distinct sequences of epithelial to mesenchymal transition (EMT) and mesenchymal to epithelial transition (MET) intermediates. The analysis yields predictions for gene expression changes consistent with curated EMT gene sets and provides a prediction of thousands of RNA transcripts through

extracellular signal-induced EMT and MET with near-continuous time resolution. The MMIST framework leverages true single-cell dynamical behavior to generate molecular-level omics inferences and is broadly applicable to other biological domains, time-lapse imaging approaches and molecular snapshot data.

# Introduction

Uncovering how cells process microenvironmental signals to activate molecular programs that lead to changes in cell state is critical for understanding mechanisms of both normal and disease physiology. Cell state is determined by molecular and cellular composition, including genome and chromatin structure<sup>1,2</sup>, proteomic<sup>3</sup> and transcript levels<sup>4</sup>, mitochondrial function<sup>5</sup>, and metabolic activity<sup>6</sup>. Cell state is intrinsically mutable and is influenced by various extracellular cues including adhesion<sup>7</sup>, mechanical signals<sup>8</sup>, soluble-ligand signaling<sup>9</sup>, and vesicle trafficking<sup>10</sup>. Here we define discrete cell states based on quantitative analysis of live-cell image data.

Single-cell omic analyses have provided an unprecedented catalog of cell states across both normal and diseased tissues<sup>11,12</sup> while spatially-resolved sequencing<sup>13</sup> and highly multiplexed imaging<sup>14–16</sup> have revealed insights into their spatial organization; however, all of these approaches lack single-cell time-ordered information, limiting the ability to draw mechanistic insights. Live cell imaging, on the other hand, readily captures cellular dynamics over timescales of seconds to days, but is limited to a small number of molecular read-outs<sup>17–20</sup>. Further, analysis of live-cell data typically relies on single timepoint “snapshots” of cell morphology or fluorescently-labeled reporters<sup>21–24</sup>. To overcome these limitations, we recently developed a morphodynamical trajectory embedding method that leverages hidden information from time-ordered live-cell trajectories, enabling improved prediction of future behavior as compared to single-snapshot-based predictions<sup>25</sup>.

It is increasingly appreciated that mechanistic understanding of both normal and diseased biological systems will require consideration of cell state dynamics. Several recent methods describe gene expression dynamics by imposing a dynamical model upon static single-cell measurements<sup>26–28</sup>, including pseudo-time estimation<sup>29,30</sup> and RNA velocity<sup>31,32</sup>. In contrast to these methods, we develop our dynamical model based upon the direct observation of single-cell dynamics obtained from live-cell imaging. Following deconvolution methods designed to estimate cell type fractions in bulk RNA-seq data, here we adopt a linear decomposition approach to identify associations between image and gene expression data<sup>33–36</sup>. The central premise of our method is that live-cell imaging and bulk molecular profiling data share commonly identifiable cell states. We now extend our morphodynamical analysis<sup>25</sup> by defining cell states based on the morphological changes of individual cells over time. To apply this definition to live-cell imaging data, we obtain quantitative dynamical models of single-cell behavior via morphodynamical trajectories. In practice, we resolve a cell state landscape over hundreds of “microstate” centers,

where transitions among microstates are described in a discrete-time Markov model framework<sup>37,38</sup>. Our data-driven modeling approach extends other statistical physics transition path-based efforts<sup>39,40</sup> by characterizing cell state changes quantitatively observed in live-cell imaging data, yielding distinct states that can be linked to unique molecular programs observed in companion molecular profiling data.

In this work, we study molecular and cellular changes in response to a panel of ligands via paired bulk RNA sequencing (RNAseq) and live-cell imaging. We focused on the well-characterized human mammary epithelial MCF10A cell line<sup>41,42</sup>, a non-transformed cell line that recapitulates key features of epithelial biology, including migration<sup>43,44</sup> and organoid formation<sup>45,46</sup>. It is also easily manipulated in a variety of assays including live-cell imaging<sup>47</sup>, knock-down<sup>42</sup>, and chemical perturbation<sup>48</sup>, and therefore is commonly used for cell biology studies. Prior studies have used MCF10A cells to probe epithelial responses to growth factors and cytokines<sup>49</sup> and to uncover molecular programs associated with EMT<sup>50–56</sup>. We explore ligand-induced cell state changes in MCF10A cells via Molecular and Morphodynamics-Integrated Single-cell Trajectories (MMIST), a novel computational methodology integrating live-cell imaging-observed dynamics and gene expression profiling. We focus on cellular response to TGFB as an illustrative example and demonstrate the quantitative linkage of EMT-associated live-cell phenotypic responses with EMT molecular programs that we validated in an external dataset<sup>57</sup>. In total, our novel data-driven modeling approach captures cell state change along sequences of cell state intermediates via live-cell and gene expression phenotypes and enables linkage of imaging and molecular data to uncover molecular correlates of distinct morphodynamic cell states.

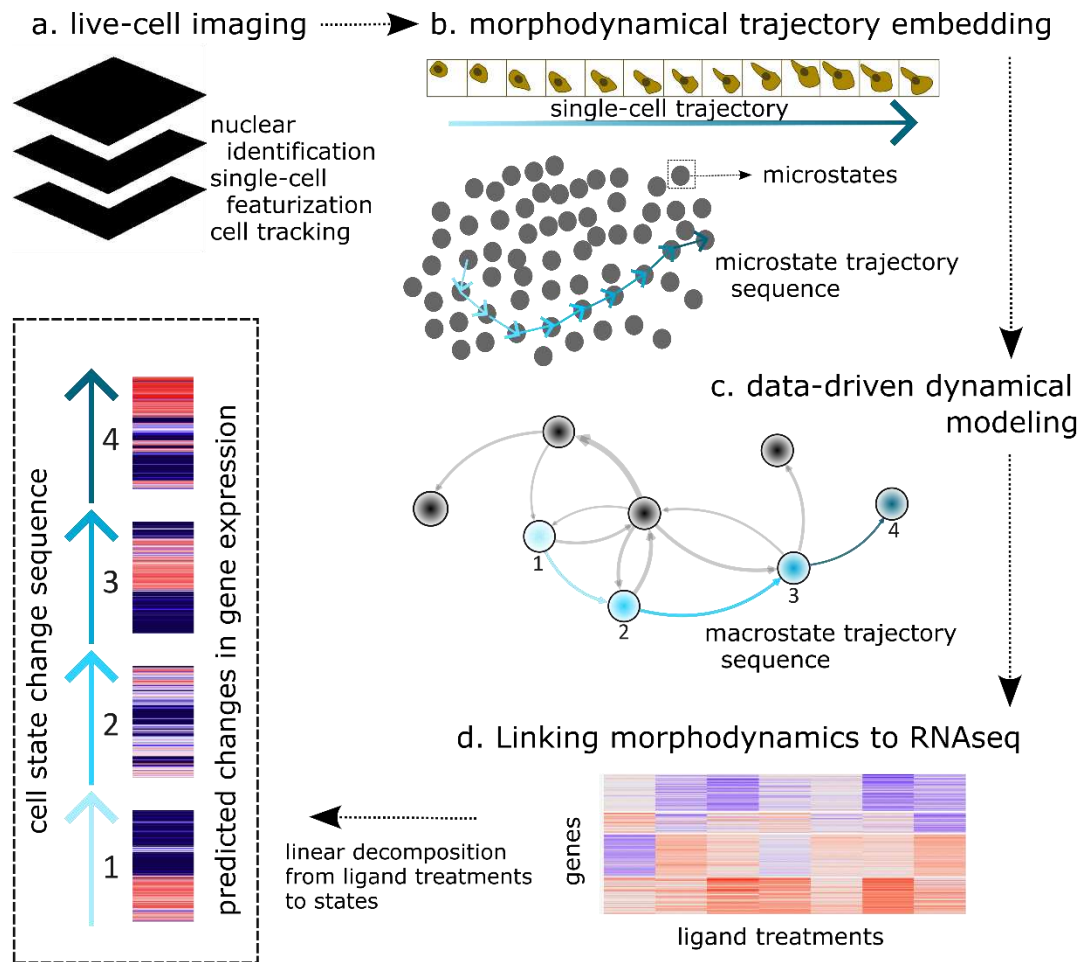
## Results

### **Experimental data to facilitate multimodal integration of morphodynamical and gene expression measurements of cell state change**

Our method is designed to infer molecular programs associated with distinct cell states by linking morphodynamic measurements acquired in live-cell imaging data to companion snapshot molecular data. We analyze a recently published LINCS MCF10A ligand perturbation dataset<sup>49</sup> which consists of paired live-cell imaging and bulk transcriptomic measurements of MCF10A cells after treatment with 6 ligands, including Epidermal Growth Factor (EGF), Transforming Growth Factor Beta (TGFB), and Oncostatin M (OSM). We also leverage live-cell image data and

transcriptomic measurements of MCF10A cells genetically engineered to express a nuclear and a cell cycle reporter<sup>58</sup> and exposed to combinations of the ligands above.

Our data analysis pipeline, illustrated in **Figure 1**, leverages companion live-cell image stacks and gene expression measured in bulk RNAseq as input, and utilizes statistical physics approaches to yield maps of cell states and their transition sequences<sup>59,60</sup>. Here we outline the major steps. (a) First, we analyze the live-cell image data to identify cell nuclei by training a virtual nuclear reporter<sup>61</sup> on paired phase contrast and nuclear reporter images, then virtually stain nuclei in the entire dataset (**Supplementary Figure 1**). We “featurize” individual cells to quantify cell shape and texture and also perform local environment featurization using Voronoi boundaries based on the nuclear centers. We track individual cells across images with Bayesian belief propagation<sup>62</sup> and compute motility as cell displacement between frames (**Supplementary Data Table 1** and **Supplementary Figure 2**). (b) Cell features are analyzed over trajectory snippets (all possible cell sub-trajectories of a particular length in a sliding window manner) utilizing our morphodynamical trajectory embedding methodology<sup>25</sup>. (c) Morphodynamical trajectories are used to build a data-driven dynamical model of cell states. (d) Cell states observed in the image data are mapped to gene transcript levels using linear decomposition. The outputs of our approach are temporal sequences of morphodynamical cell state changes and their associated gene expression levels.



**Figure 1: MMIST approach to link live-cell imaging to molecular read-outs.**

a) Live-cell imaging of MCF10A cells after treatment with a panel of microenvironmental ligands. Nuclei are identified using a convolutional neural network, and single-cells are featurized and tracked through time. b) Single-cell features are concatenated along single-cell trajectories to construct the morphodynamical trajectory embedding. c) Dynamical models learn cell states and cell state change sequences in the morphodynamical landscape. d) Cell state populations are used as a linear decomposition of bulk gene expression measurements to predict the gene expression programs underlying cell state change.

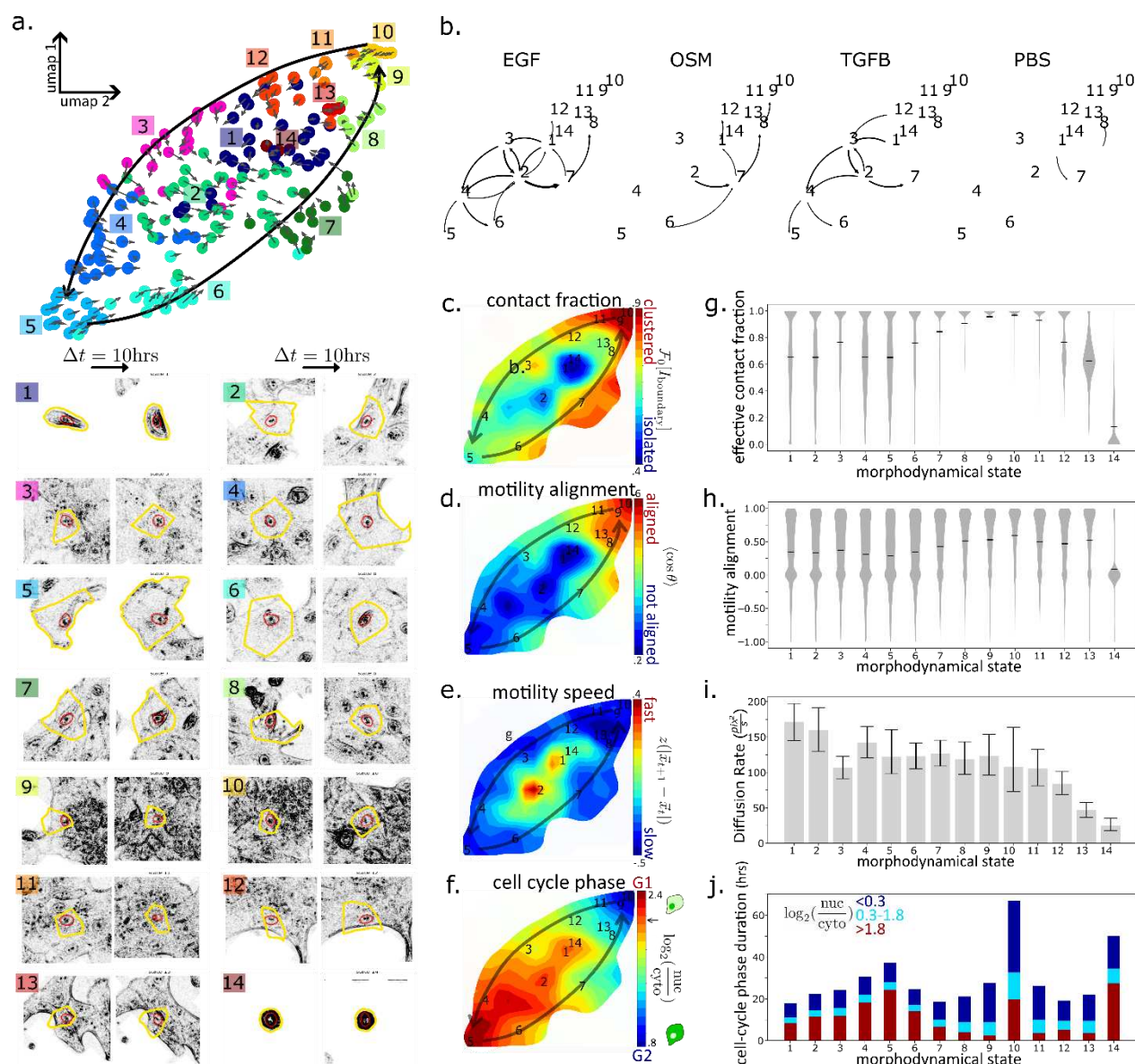
### Single-cell trajectories define morphodynamical cell states

Our goal is to group cells into states with shared dynamical progression—i.e., those that cluster together based on a similar time progression of shape, texture, and motility features. The morphodynamical trajectory space is a time-concatenation of image-based features<sup>25</sup> in which we place hundreds of “microstate” centers via clustering. We then count transitions among microstates to build a data-driven transition matrix Markov model of cell state progression<sup>37,38</sup>. Next, microstates are grouped into coarser “macrostates” using a spectral clustering

procedure<sup>63,64</sup>. We refer to these macrostates as morphodynamic states, or simply states. The eigenfunctions of the dynamical model represent dynamical motifs, which we visualize using UMAP<sup>65</sup> dimensionality reduction to facilitate interpretation of cell states (**Figure 2a**). Ligand treatments induce unique cell state changes and transition flows as compared to negative control (PBS) (**Figure 2b**). The complete set of ligand-dependent cell state populations are shown in **Supplementary Figure 4c**, and population distributions and transition flows are shown in **Supplementary Figure 5**.

The derived states largely resolve differences in morphodynamical properties such as the cell-cell contact fraction, local alignment of cell-cell motility, motility speed, and cell-cycle phase (**Figure 2c-j**). Cell states 5 and 10 represent two distinct morphodynamic states bracketing the morphodynamical cell state space. State 5 is characterized by mesenchymal-associated features such as lower local alignment of cell motility, more extended cytoskeletal features, greater cell spreading, and an extended G1 cell cycle duration (**Figure 2j**); this state population increases under TGFB containing treatments. In contrast, state 10 is characterized by many epithelial-associated features, including increased multicellular clustering and collective motility which are increased after treatments that include OSM; these represent an altered epithelial phenotype that maintains epithelial-associated characteristics.<sup>49</sup> Between these two states, we observe intermediate states with short cell cycle duration (**Figure 2f,j**), increased motility (**Figure 2e,i**) and the fewest cell-cell contacts (**Figure 2c,g**). Under EGF treatment, which is typically added to MCF10A cell culture medium<sup>41</sup>, cells transition between these intermediate states (**Figure 2b**). Thus, based upon morphodynamical features, the derived cell state space matches the well-described biological framework of epithelial and mesenchymal cell states<sup>66</sup>, including extended G1 duration in the mesenchymal state<sup>67–70</sup>.





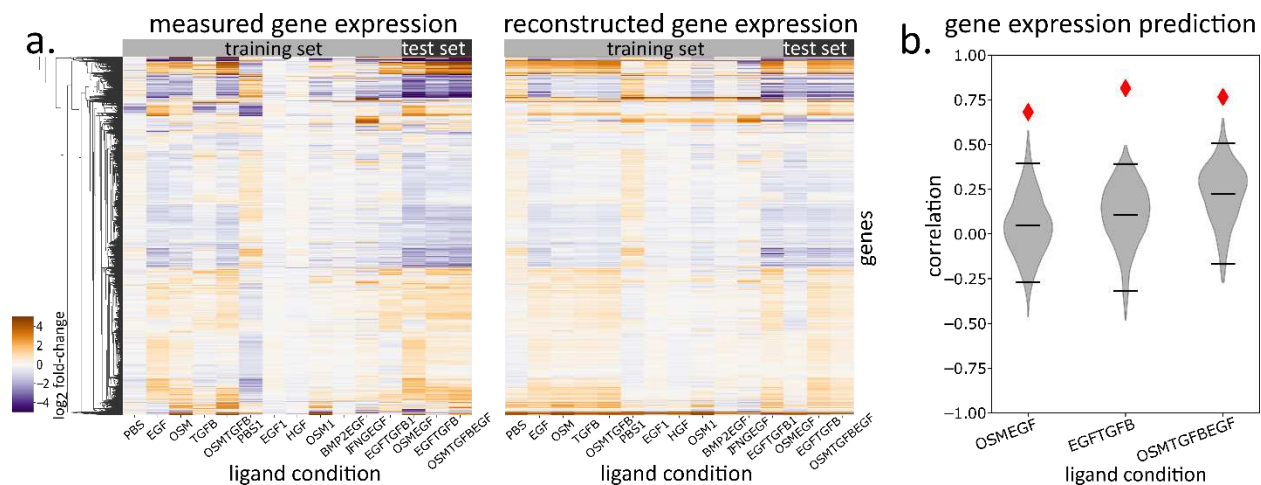


## **Morphodynamic trajectories reveal transcriptional dynamics via state mapping across modalities**

Motivated by the observation that the morphodynamic cell states recapitulate aspects of EMT, we next sought to identify the underlying molecular programs associated with cell state transitions. The process relies on having both morphodynamical observations and molecular measurements for an identical set of experimental treatments to enable linkage of RNA transcript levels to the morphodynamical states delineated above. The primary assumption – which can be considered a hypothesis being tested – is that an observed morphodynamical state corresponds to the same RNA levels regardless of the ligand treatment. Consider the example of linking RNA levels to two distinct states (motile and non-motile), where the cell state frequencies are modulated by ligand treatment. If ligand A induces an increase in the motile cell state population as compared to B and *also* higher RNA levels for gene X, then we infer that motility is linked to expression of gene X. This qualitative idea can be made exact in a simple linear algebra framework by decomposing each measured average transcript level as a linear sum over morphodynamical state populations (**Supplementary Figure 4c**) and gene expression profiles.

We first validate the linear population matching approach by assessing its capability to predict withheld gene expression levels in ligand combination conditions. The method requires at least as many paired live-cell imaging and RNAseq measurements as states, so we performed a separate clustering into 10 morphodynamical cell states, allowing us to withhold the OSM+EGF, EGF+TGFB, and the triple combination OSM+EGF+TGFB RNAseq data from the training set used to extract morphodynamical cell state gene expression profiles. The morphodynamical cell state populations from the live-cell imaging in the withheld test set treatments, combined with morphodynamical cell state decomposed gene expression levels from the training set, enable a prediction of the RNAseq in the test set (**Figure 3a**). The predictive capability of the model is maximized at a trajectory length of 10 hours, where these predictions yield a Pearson correlation  $>0.7$  to the test set gene expression and high significance compared to a null model with random state populations (p-value $<0.001$ , upper-tailed test, **Figure 3b**); correlation coefficients and true positive rates to predict up/down regulation for different trajectory lengths are in **Supplementary Data Table 2**. Performance exceeding the random null model demonstrates that the defined states exhibit treatment-independent association with the inferred expression levels. These

findings provide support for the validity of our approach to link morphodynamical states observed in image data to companion molecular measurements.



**Figure 3: Morphodynamical cell states predict ligand combination gene expression.**

a.) Validation of model gene expression predictions: measured and model-reconstructed gene expression at 24hrs for every experimental condition, including training set (light gray) and test set conditions, b.) Correlation between measured and model-predicted gene expression (red diamonds), and null estimates using random state populations (gray violin plots).

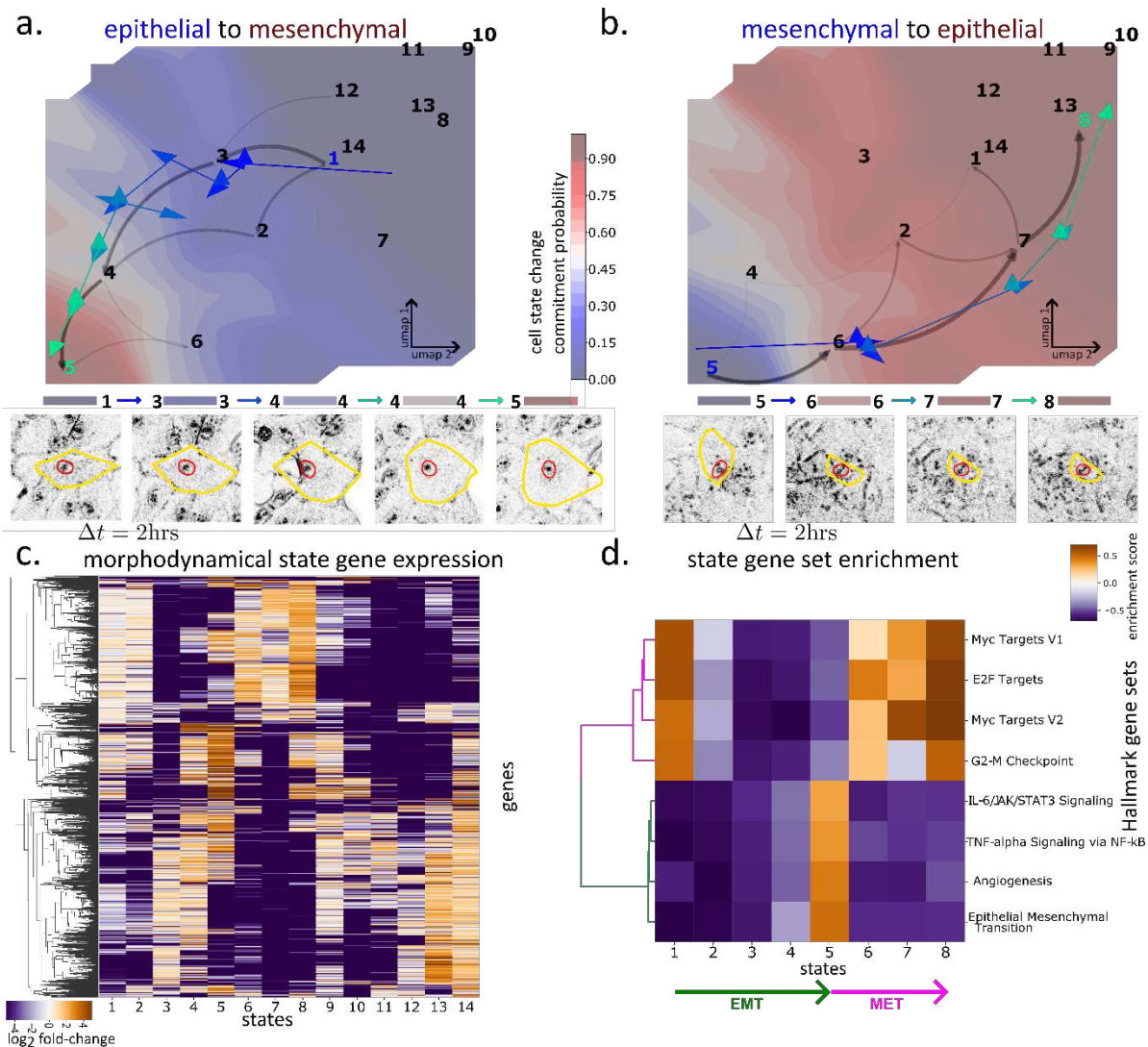
## MMIST identifies ligand-induced EMT and MET morphodynamical cell state change sequences

We next focused on epithelial-mesenchymal morphological features as an illustrative use case of MMIST. Analysis of the morphodynamical cell states revealed features associated with canonical epithelial and mesenchymal states, including changes in cell-cell motility alignment and cell clustering (**Figure 2c-j**). These findings are consistent with the observation that microenvironmental signals can strongly modulate differentiation state of MCF10A cells—for example to form epithelial-differentiated multicellular acinar structures in 3D cultures<sup>71</sup> or to promote a mesenchymal phenotype under TGFB treatment<sup>52,53,56,57</sup>. We used our framework to examine the relationship between these states and the influence of microenvironmental signals in mediating transitions between them, which builds on prior studies of epithelial-mesenchymal transition (EMT) and mesenchymal-epithelial transition (MET)<sup>26,40,52,57</sup>. To study EMT in our framework, we assigned state 1 as the most highly populated state at the initial trajectory time window (10 hours), while state 5 was assigned as mesenchymal due to its morphological features and enrichment in the TGFB condition. We set the most highly populated state at the initial time window as the initial state to facilitate identification of the most common ligand-induced state

transitions ending in the mesenchymal-like state 5. For MET, we assign state 8 as the final state because it is enriched over time in the control EGF treatment as cells reach confluence (**Supplementary Figure 6**).

We observe robust but distinct state transition sequences for the EMT and MET transitions, consistent with a highly driven nonequilibrium system<sup>72,73</sup>. For EMT the dominant sequence of states is (1,2,3,4,5) but for MET, the primary sequence is (5,6,7,8), shown in **Figure 4a,b**. Transitions back to state 1 are common in most treatments (**Figure 2b** and **Supplementary Figure 5**). The dominant sequences of state changes are robust across ligand treatments, though the probability of specific state-to-state transitions varies. For instance, OSM treatment drives most cells towards dense and collectively migrating epithelial-like clusters (state 10), but for the rare cells which do reach state 5 from state 1, the dominant sequence of states remains the same (**Supplementary Figure 7**).

MMIST revealed unique expression patterns associated with each morphodynamical cell state (**Figure 4c**). We performed gene set enrichment over the Hallmark gene sets<sup>74</sup> on the derived morphodynamical state gene expression profiles. The morphodynamical state-decomposed gene expression along the EMT state change sequence shows a transition from a proliferative program enriched for Hallmark Myc Targets V1 and V2, E2F targets, and G2-M transition, to a mesenchymal program enriched for IL4/JAK/STAT3, TNFA via NFKB, Angiogenesis, and Epithelial to Mesenchymal Transition (**Figure 4d**). This switch from a proliferative program to a mesenchymal gene expression program augments our observation that cell-cycle phase durations co-varied with mesenchymal-like features observed in the live-cell data (**Figure 2j**).



**Figure 4: Morphodynamical cell states predict ligand combination gene expression.**

a.-b.) Cell state change pathways (black arrows; thickness proportional to probability flux carried by each state-to-state transition) based on cell states from Figure 2a, and cell state change commitment probability (blue to red) in EGF (reference positive control) condition. Also shown are representative single-cell trajectory (dark blue to turquoise arrows, 30min timestep) and cell images (1 hr between images). c.) Differential gene expression in each morphodynamical cell state (top 8000 most variable genes), with magenta and green labels corresponding to assignment to Hallmark gene sets labeled in d., and transcription factors labeled on y-axis. d.) Hallmark gene set enrichment over EMT/MET associated cell states.

## Near-continuous gene expression time evolution prediction during TGFB-driven EMT

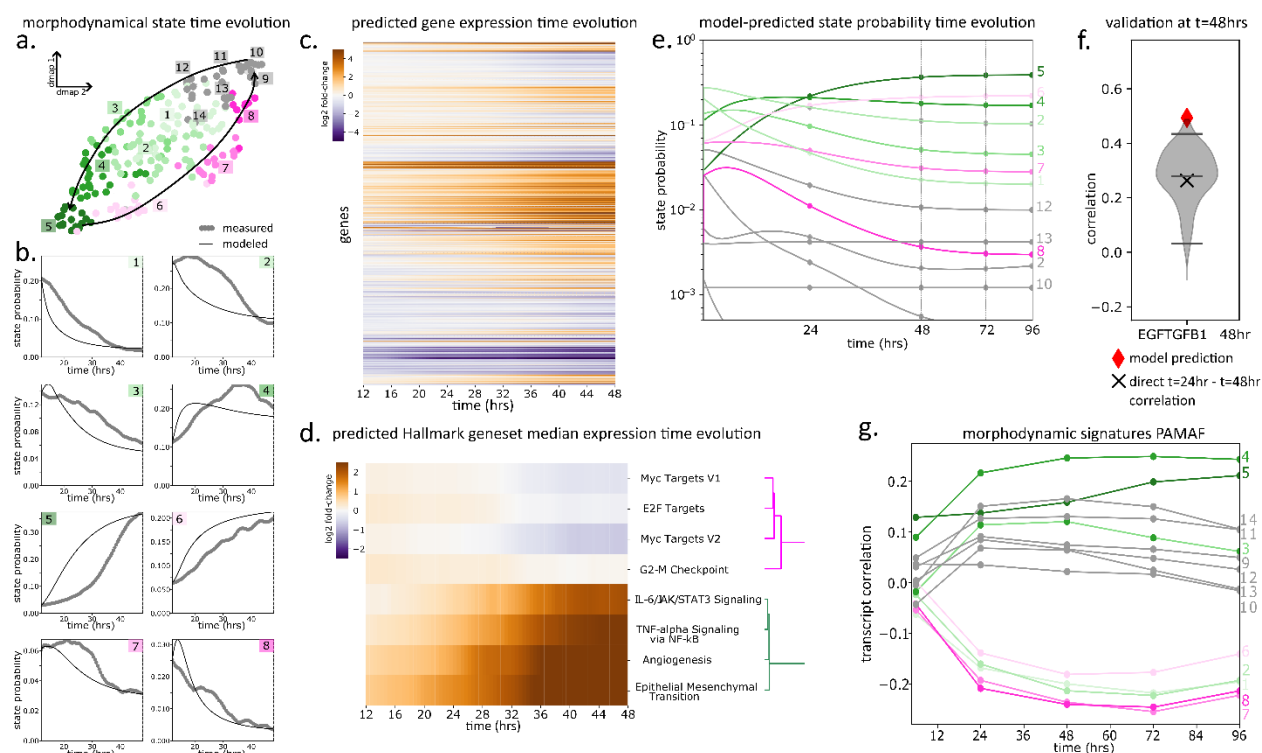
MMIST yields near-continuous evolution of morphodynamical cell state populations by counting transitions between microstates extracted from single-cell trajectories. For example, EGF+TGFB leads to an increase in mesenchymal-like states 4, 5, and 6, whereas these states are decreased after EGF-only treatment (**Figure 5a**). A key component of MMIST is to model cell state evolution with a Markov model. We assessed this aspect of the model by comparing inferred and modeled state populations as a function of time after EGF+TGFB treatment. The model largely reproduces the morphodynamical state population trends observed in the live-cell imaging experiment, supporting the validity of our Markov assumptions (**Figure 5b**).

Our computational framework enables a prediction of gene transcript levels at the same near-continuous time intervals as those measured in the live-cell image data. Conceptually, we do this by leveraging the observation that each morphodynamical state is associated with a gene expression profile and then predict the bulk gene expression over time by computing a weighted sum of the states observed in each treatment condition (**Figure 5c**). Under EGF+TGFB, our model predicts a continuous shift in multiple gene programs, including decreases in proliferation-associated programs and increases in mesenchymal-associated programs (**Figure 5d**). MMIST can also be used to predict future, unmeasured shifts in cell state populations. For example, the model predicts large shifts in state populations between 0-48H, which we observed experimentally; however, it also predicts continued subtle shifts in state populations beyond the 48H duration of the experiment (**Figure 5e**). We next assessed the ability of our model to predict unseen changes in gene expression programs. Here, we trained our model with RNAseq data collected at 24H post-treatment, then used it to predict gene expression profiles at 48H based on the predicted morphodynamical state populations shown in **Figure 5e**. We assessed our predictions by computing the correlation between experimentally measured and predicted expression profiles, after normalizing to t=0H. The correlation between predicted t=48H gene expression profile and the withheld t=48H RNAseq data is ~0.5 (**Figure 5f** and **Supplementary Figure 4**). In contrast, t=24H and t=48H experimentally measured RNAseq profiles show correlations of ~0.25, indicating that MMIST predictions can capture molecular programs associated with morphodynamic state change.

The transcriptional programs associated with TGFB-driven EMT have been previously investigated in MCF10A cells, and datasets generated through these efforts provide a useful tool for independent validation of our model<sup>52,53,56</sup>. To evaluate the EMT-associated signature

extracted via our morphodynamical analysis, we compare our results to a recently published, independent, time-resolved gene expression dataset of MCF10A cells treated with EGF+TGFB then harvested for molecular profiling at multiple timepoints including 24, 48, 72 and 96 hours post-treatment, (“PAMAF” data)<sup>57</sup>; this dataset lacked companion live-cell image data. We first assess the biological significance of the model-assigned morphodynamic states based on gene expression levels, finding positive correlation between PAMAF measurements and mesenchymal morphodynamical cell states 4 and 5 after EGF+TGFB treatment (**Figure 5g**). Consistent with this, epithelial states 6,7, and 8 are among the least correlated. Together, these findings provide support for the robustness of MMIST to identify meaningful biological signals that can be validated in independent data sets.





**Figure 5: Morphodynamical model predicts EGF+TGFB-induced EMT gene expression time evolution**

a) Morphodynamical states, which are numbered 1-12 and color-coded (mesenchymal: green, epithelial: purple). Color labels for the states are consistent throughout figure. b) State probability time evolution, measured (grey dots) and model-derived (black lines). c) Prediction of gene expression over time at 30-minute intervals using morphodynamical state prediction and live-cell imaging measured state probabilities, showing the top 8000 most variable genes (top) and d) summarized to Hallmark gene sets (bottom). e) Model-predicted state probability time evolution over 96 hours, trained from live-cell imaging over 48 hours. f) Correlation between measured and model-predicted gene expression at t=48H (red diamond) based on training data from t=24H, relative to null models with random state probabilities (gray distribution). Also shown: correlation between t=24H and t=48H gene expression (black X). g) Correlation between predicted morphodynamical state gene signatures and PAMAF measurements out to 4 days.

## Discussion and Conclusion

Single-cell sequencing and spatial omics methods have provided detailed molecular profiling of cellular heterogeneity in single time-point snapshots<sup>13</sup>. However, there are no methods yet that yield time-resolved molecular profiles with a similar level of detail. RNA velocity and other algorithmic methods attempt to infer dynamics from fixed measurements<sup>26–31</sup>, but they lack a



direct mapping to observed single-cell dynamics. Here we have presented a step in the direction of linking live-cell dynamics to deep molecular profiling, capturing sequences of morphodynamical cell state changes mapped to comprehensive gene expression profiles. Our approach provides a direct map from live-cell derived single-cell dynamics to gene expression for a small number of morphodynamical cell states defined through assessment of perturbation responses.

We utilized a paradigm of cell behavior in which individual cells transition between different morphodynamic states with treatment-specific dynamics and state frequencies. Thus, we employ a trajectory space that is common to all observed experimental treatments, where ligand perturbation alters the rates of cell state changes. This report demonstrates the value of this paradigm, as it enables mapping of complex, spatiotemporal phenotypes to gene transcript levels. One limitation of the present model is that it is restricted to the range of behaviors observed for the particular cell type (MCF10A) under the treatments examined and does not represent a comprehensive assessment of all possible cell states. Thus, the derived (coarse-grained) dynamical models are incomplete. As live-cell information increases, for instance via the incorporation of multiplexed live-cell reporters and deep-learning based image featurization<sup>75–78</sup>, integration with fixed single-cell and spatial omics profiling at endpoint may require a separation of shared information across cell populations from unique information to each single-cell<sup>79</sup>.

From a physical theoretical point of view, the transition “mechanism” of a dynamical process is defined via the set of trajectories connecting two states of interest<sup>59,80–82</sup>, for instance epithelial and mesenchymal cell states. The single-cell trajectory set that connects these basins contains the set of intermediate transition states<sup>39</sup>. Here, we have captured sequences of EMT and MET intermediates, consistent with the emerging view of epithelial and mesenchymal states as a continuum<sup>11,26</sup>. It is an open question of whether characterization of transition intermediates will yield insight into cell state control, which could inform the control of EMT-driven processes during development or disease progression, such as tumor invasion<sup>83,84</sup>. Future studies could extend our findings by employing inhibitor or gene knockout approaches to functionally assess EMT transition intermediates predicted to be critical for cell state control.

Cell state biomarkers can predict sensitivity to targeted drugs<sup>85,86</sup>, and are expressed in a spatially organized manner in both healthy and diseased tissues<sup>87,88</sup>. Morphodynamical cell state definitions can expand upon known biomarker-based cell states, providing a prediction of the dynamical responses to biological manipulation. We expect that the linking of morphodynamics

to gene expression changes, in spatial context, will lead to a deeper understanding and control of cell state change in complex tissue and tissue-like environments.

Characterization of the transition mechanism via live-cell image-based trajectories, such as we have presented, is not a mechanistic explanation at the molecular level. Time-ordered single-cell trajectories of the quantity of molecular species, such as gene transcripts, imply but do not prove causality. We speculate that utilizing molecularly detailed single-cell trajectory data to constrain mechanistic models could provide prediction of causal molecular relationships that could be experimentally validated. Our data-driven approach, as presented, does not yield a prediction for unmeasured perturbations, for instance response to different ligands or drugs. We speculate that mechanistic models<sup>89–92</sup>, trained using the type of detailed trajectory data at the molecular level we have presented here, may enable prediction of cell behavior in unseen contexts.

Live-cell phenotypic response to ligand perturbation is well-described by our single-cell morphodynamical trajectory-based data-driven modeling approach and enabled a mapping between live-cell phenotype and time-dependent gene expression changes. Our models yielded a validated prediction of near-continuous gene expression levels during ligand-driven EMT/MET in MCF10A cell culture. MMIST can be applied generally to characterize cell state changes in fundamental biology and, potentially, in various disease settings.

## Methods

**MCF10A Cell Culture** MCF10A cells were cultured in growth media composed of DMEM/F12 (Invitrogen #11330-032), 5% horse serum (Sigma #H1138), 20 ng/ml EGF (R&D Systems #236-EG), 10 µg/ml insulin (Sigma #I9278), 100 ng/ml cholera toxin (Sigma #C8052), 0.5 µg/ml hydrocortisone (Sigma #H-4001), and 1% Pen/Strep (Invitrogen #15070-063). For all ligand response experiments, cells were seeded in growth media in collagen-coated well plates and allowed to attach for 6-hours. Cells were then washed with PBS, and growth media was replaced with growth-factor free media lacking EGF and insulin. After an 18-hour incubation, cells were treated with ligands in fresh growth-factor free media. Seven different ligand conditions were tested at concentrations previously determined to elicit maximal cell responses<sup>49</sup> (EGF 10 ng/ml (R&D Systems #236-EG), OSM 10 ng/ml (R&D Systems #8475-OM), TGFB 10 ng/ml (R&D Systems #240-B), EGF 10 ng/ml + OSM 10 ng/ml, TGFB 10 ng/ml + EGF 10 ng/ml, OSM 10 ng/ml + TGFB 10 ng/ml, TGFB 10 ng/ml + EGF 10 ng/ml + OSM 10 ng/ml). Wild type MCF10A cells were a generous gift from the Gordon B. Mills lab, and were used for all RNA-seq

experiments. For live-cell imaging experiments, parental WT MCF10A cells were genetically modified as described below.

**Live-cell imaging** To assess cell-cycle responses to ligand treatments, MCF10A cells were genetically modified to stably express the HDHB cell-cycle reporter<sup>58</sup> and a red nuclear reporter. The methodology used to generate the reporter cell line has been described previously<sup>93</sup>. Reporter cells treated with ligand were imaged every 15 minutes for 48 hours with an Incucyte S3 microscope (1020x1280, 1.49  $\mu\text{m}$ /pixel). Three channels were collected -- phase contrast, red (nuclear) and green (cell-cycle) -- for four fields of view per well. The initial frame coincided with the addition of the ligands and fresh imaging media. A previously published dataset of live-cell imaging results (imaged every 30 minutes for 48 hours) was also analyzed in this study, specified here by appending a 1 to the treatment condition (e.g. EGF1)<sup>49</sup>. All matching ligand treatments utilized identical ligand sources and concentrations in both datasets. This additional dataset was generated from WT MCF10A cells dosed with a broad panel of single ligand treatments, using similar cell culture and imaging techniques. Further experimental protocols from this study can be found in detail at the publicly available Synapse database<sup>49</sup>.

**RNAseq** Detailed description of sample preparation, processing, and alignment can be found in Gross et al<sup>49</sup>. For each ligand treatment, we performed a differential expression analysis from time zero controls on the RNAseq gene-level summaries with the R package DESeq2 (1.24.0), with shrunken log2 fold change estimates calculated using the apeglm method. We applied a minimum expression filter such that  $\log_2(\text{TPM}) > 0.5$  in at least 3 measurements over treatments and replicates (with TPM transcripts per million), yielding 13,516 genes with measured differential expression from control used in our analysis.

**Image preprocessing** Foreground (cells) and background pixel classification was performed using manually trained random forest classifiers using the ilastik software<sup>94</sup>. Images were z-normalized (mean subtracted and normalized by standard deviation). In cell images, absolute values of these z-normalized pixel values are shown (white to black). Image stacks were registered translationally using the pystackreg implementation of subpixel registration<sup>95</sup>.

**Nuclear segmentation** A convolutional neural network was trained to predict the nuclear reporter intensity from the matched phase contrast images for imaging data of WT MCF10A cells with no nuclear reporter. In the EGF, OSM, and TGFB conditions, 4 image stacks (12 total) were used to

train the FNET 3D reporter prediction CNN from the Allen Cell Science Institute<sup>61</sup>, with time as the third dimension rather than z-dimension. This trained CNN was then used to predict nuclear reporter channel from the bright-field image over all image stacks in datasets. See Supplementary Figure 1 for representative nuclear reporter prediction and comparison to ground truth. Nuclear segmentations were generated by performing a local thresholding of the image within 51 pixel-sized windows at 1 standard deviation of intensity. Segmentations were filtered for a minimum size of 25 pixels and a maximum size of 900 pixels, see Supplementary Table 1 for segmentation performance. To capture features including the local environment around a single nucleus, the image was partitioned into Voronoi cells around each nuclear center, with background classified pixels removed. Image preprocessing and segmentation scripts can be found on the github repository, see data and code availability.

**Cell featurization** Single-cell featurization was performed on the Voronoi-partitioning of the image by nuclear center. Cell features are described in detail in Copperman et al.<sup>25</sup> and repeated here for convenience. Morphology features were obtained as follows: segmented cells were extracted, and mask-centered into zero-padded equal sized arrays larger than the linear dimension of the biggest cell (in each treatment). Principal components of each cell were aligned, and then single-cell features were calculated. Zernike moments (49 features) and Haralick texture features (13 features) were calculated in the Mahotas<sup>96</sup> image analysis package. The sum average Haralick texture feature was discarded due to normalization concerns. Rotation-invariant shape features (15 features) were calculated as the absolute value of the Fourier transform of the distance to the boundary as a function of the radial angle around cell center<sup>97</sup>, with the set of shape features normalized to 1. The cell environment was featurized in a related fashion. First, an indicator function was assigned to the cell boundary with value 0 if the boundary was in contact with the background mask, and value 1 if in contact with the cell foreground mask. The absolute value of the Fourier transform of this indicator as a function of radial angle around cell-center then featurized the local cell environment (15 features), with the sum of cell environment features normalized to 1. Note the first component of the cell environment features is practically the fraction of the cell boundary in cell-cell contact. The high-dimensional cell feature space was dimensionally reduced using principal component analysis (PCA), retaining the largest 11 eigen-components of the feature covariance matrix (spanning all treatments and image stacks) which captured >99% of the variability.

**Motility features** Cell motility was characterized in a single-cell manner, referenced both to the image frame and relative to neighboring cells. Single-cell displacement  $\vec{\Delta x}$  between tracked frames was z-normalized, and cells which could not be tracked backward for a frame had unrecorded displacements and were not used in our analysis. The local motility alignment of a single-cell to the local neighborhood of contacting cells (sharing a Voronoi boundary) was measured by extracting the cosine of the angle between the single-cell and direct neighbors via  $\hat{p}_1 \cdot \hat{p}_2$  with  $\hat{p} = \vec{\Delta x} / |\vec{\Delta x}|$ . Local contact inhibition of locomotion was measured via the higher-order vector formed by  $(\hat{p}_1 - \hat{p}_2) \cdot \hat{r}_{12}$  with  $\hat{r}_{12}$  the separation vector between cells<sup>98</sup>. Neighborhood averages were taken via the Voronoi partition, averaged over neighbors and weighted by the relative length of the boundary to each neighbor, see Supplementary Figure 3.

**Batch normalization** Single-cell featurization can depend in subtle ways upon the imaging treatment and sample batch. To normalize these effects we utilized a batch normalization procedure at the single-cell feature level. For each morphology feature, we utilized a histogram matching procedure between negative control (PBS) treatments. We then fit a linear model to the histogram-matched distributions, and applied this linear model between sample batches, see Supplementary Figure 8.

**Cell tracking** To follow single-cells through time to extract the set of single-cell trajectories for morphodynamical trajectory embedding, we utilized a Bayesian likelihood-based approach implemented in the btrack software package<sup>62</sup> using default parameters. This Bayesian approach was applied for each frame over a 12 frame window, and then successful tracks over each pair of successive frames were extracted. See Supplementary Table 1 for manual validation of tracking performance.

**Morphodynamical trajectory embedding** To maximize the single-cell information, we extended single-timepoint morphology and motility features over single-cell trajectories using a delay-embedding approach, described in Copperman et al.<sup>25</sup> In brief, single-cell features including motility features, but excluding cell-cycle features, were concatenated along the trajectory length to form morphodynamical feature trajectories. We tested multiple trajectory lengths and selected a trajectory length of 10 hours where the best prediction of withheld treatment combination RNAseq was obtained, see Supplementary Table 2. We utilized a dynamical embedding approach described below to cluster trajectories and visualize this space, and did not perform any

further dimensionality reduction upon the trajectory concatenated morphological feature PCAs and motility feature trajectories prior to dynamical model building.

**Data-driven dynamical Markov state model** To capture dynamical properties within the morphodynamical space, we constructed a transition matrix Markov model within the trajectory embedding space. The embedded space was binned into “microbins” using k-means clustering with  $k = 200$  clusters. Results using 50, 100, 200, and 400 clusters are qualitatively similar. In this discrete space, a transition matrix  $T$  between bins was estimated from the set of transition counts  $C_{ij}$  from microbin  $i$  to  $j$  as  $T_{ij} = C_{ij}/C_i$  with  $C_i = \sum_j C_{ij}$ . This accounting was agnostic to cell birth and death processes, yet we observe our model well reproduces morphodynamical state evolution, see Figure 5b.

**Dynamical features** To evaluate live-cell behavior via characterization of shared dynamics, we have applied a dynamical featurization approach via the data-driven transition matrix model. Using a transition matrix model constructed from all possible single-cell trajectory steps in the in the microbinned trajectory feature space, we construct the Hermitian extension  $H = \frac{1}{2}[(T + T') + i(T - T')]$  with  $T'$  the transpose of the transition matrix  $T$ , this approach numerically stabilizes the eigendecomposition and provides all real eigenvalues for unambiguous ordering of eigencomponents<sup>99</sup>. We retain 15 dominant eigencomponents (see Supplementary Figure 10), and concatenate real and imaginary parts of eigenvectors to construct a 30-dimensional characterization of each microbin center. To visualize the dynamical trajectory space, we apply UMAP dimensionality reduction of the microstate eigenvector components to 2 components. Average flows in the UMAP space are calculated via calculating microstate dependent average displacements via the transition matrix  $\langle x_i \rangle = \sum_j (x_i - x_j)T_{ij}$  and averaging over 10 nearest microstate neighbors for smoothness. We note that UMAP flows were used only for visualization, not featurization.

**Morphodynamical cell states** As a tool for reducing complexity and extracting biological meaning in the morphodynamical embedding space, we defined a set of macrostates by clustering together microbins using dynamical similarity. We utilize the eigencomponents of  $H$  (Hermitian extension of the transition matrix  $T$ , see Dynamical Embedding) and perform k-means in the kinetic motifs. We utilize a lower cutoff of 0.015 for the total fraction of cell trajectories assigned to each state; if a microstate has too few trajectories assigned, then it is combined with



its nearest neighbor by Euclidean distance in the space of dynamical motifs. k-means clusters are increased until the requested number of states with minimum fraction assignment is obtained. We then evaluated the capability of the derived macrostates to describe the state-change dynamics by evaluating the sum of timescales captured in the microstate transition matrix model, related to the VAMP score<sup>100</sup>. We observe a rapid increase in score increasing to 10 states and continued increase beyond 15 states, see Supplementary Figure 9. Note that the macrostates, like the features themselves, were not designed or optimized for the task of predicting RNA levels.

**Cell-state change pathways** To extract the sequences of morphodynamical cell states under EMT/MET, we adopted a transition path approach to calculate committers and state change sequences utilizing our data-driven Markov model<sup>59</sup>. Transition matrices were constructed between morphodynamical cell states (macrostates), and flux analysis was carried out using the PyEMMA analysis package<sup>60</sup>; all pathways carrying flux between sets of initial and final states were evaluated to find dominant state change sequences. Commitor probabilities (for reaching the final state before returning to the initial state) were highly dependent upon culture treatment, but cell state change sequences were quite robust to culture treatment, see Supplementary Figure 7.

**Cell-cycle reporter analysis and dynamical modeling** To capture cell-cycle dynamics from the HDHB reporter images, we adopted a similar data-driven modeling approach as we took in defining the morphodynamical cell states. Reporter levels in the nuclear and cytoplasmic compartments were extracted, and the ratio of these reporter levels was used as a self-normalizing readout of cell-cycle state, where exclusion of HDHB from the nucleus is known to correlate with G2 cell-cycle state, with maximal nuclear correlation occurring abruptly at mitosis and decreasing gradually from G1 to S, and with minimal nuclear signal at G2<sup>93</sup>. To divide reporter ratio values into cell-cycle stages, we utilized our Markov state modeling and dynamical embedding procedure, first building a microbin model with 50 bins evenly spaced throughout the range of reporter ratio values, then dividing these into 4 macrostates via k-means clustering in dynamical motifs, see Supplementary Figure 10. We then calculated mean first passage times using PyEMMA between cell-cycle stages as a readout of cell-cycle stage lifetimes in each of the morphodynamical cell states.

**Bulk RNAseq reconstruction** To capture the biological drivers of morphodynamical cell state changes, we mapped our morphodynamical cell states to RNAseq-based gene expression



profiles. We adopted a linear decomposition approach. If cells in treatment A are subdivided into a set of states  $s$  with known state populations  $p_s^A$  such that  $\sum_s p_s^A = 1$ , and the state and treatment dependent gene levels are known, a bulk measurement of the  $i$ th gene can be reconstructed exactly as  $\langle g_i^A \rangle = \sum_s p_s^A g_i^{s,A}$ . We approximate this exact expression by making the assumption that all cells in state  $s$  under each treatment have identical gene expression, i.e., that  $g_i^{s,A} = g_i^s$  regardless of A, for every  $s$ . The utility of this approximation can be evaluated via our results, and is equivalent to letting the states form a non-negative matrix factorization of the bulk expression. Under this assumption, we have a linear system of equations connecting state populations and state gene expression levels  $\{\langle g_i^A \rangle = \sum_s p_s^A g_i^s, \langle g_i^B \rangle = \sum_s p_s^B g_i^s, \dots\}$ , one equation for each treatment A,B,C,... based on treatment-specific cell state populations  $p_s^A, p_s^B, \dots$  directly measured via live-cell imaging and morphodynamical analysis, and with paired bulk RNAseq measurements  $\langle g_i^A \rangle$ . If there are as many measurements as states, this linear equation can be inverted for the gene expression profiles in each state,  $g_i^s$ . If there are less states than treatments and the solution is over-determined, we obtain the solution over all possible combinations of treatments and average over the results. In practice, true solution of the linear system would yield negative gene levels, so we do a least squares minimization with the constraint of positive gene levels. We use fold-changes rather than absolute gene levels to preserve the batch and replicate normalization, this normalization does not affect the system of equations as it enters on both sides of the equality. To validate our state decomposition of measured bulk RNAseq pipeline, we split our data into training sets and validation sets. State gene expression levels are trained from the training set gene levels only, and gene expression for withheld test set conditions are then predicted via the measured morphodynamical cell state populations. Null model predictions are constructed from random state populations combined with previously estimated state-specific gene levels (from true populations) as a measure of how unique the measured state populations are at predicting the test set gene expression.

**Gene set enrichment** To interpret morphodynamical cell state gene expression profiles, we performed gene set enrichment analysis via the pyGSEA package<sup>101</sup>. We utilized the preranked algorithm, sorting genes via the predicted gene expression levels in each morphodynamical cell state. We ran gene set enrichment using the Hallmark gene sets<sup>74</sup>, which broadly capture well-studied biological processes and cell signaling activity.

## **Acknowledgements**

We thank Joe Gray for contributions to project conceptualization, Mark Dane for technical guidance and assistance accessing LINCS data, David Aristoff and Gideon Simpson for mathematically oriented discussions, and John Russo and Luke Ternes for input regarding computational implementation. J.C. was supported by the Damon Runyon Cancer Research Foundation Quantitative Biology Fellowship DRQ-09-20. Y.H.C is supported in part by the National Cancer Institute (U54CA209988, U2CCA233280, U01 CA224012). D.M.Z. acknowledges support from the National Science Foundation (MCB 2119837). These studies were supported in part by NIH research grants U54-CA209988 and U54-HG008100, and the Anna Fuller Foundation.

## **Competing Interest**

The authors declare no competing interests.

## **Data and Code Availability**

All codes and scripts to perform the analysis in this work can be found at the project github repository.

<https://github.com/jcopperm/celltraj>

LINCS MCF10A Molecular Deep Dive data is available in some formats from the synapse database<sup>49</sup> and additional data is available upon request.

## References

1. Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome. *Nature* **489**, (2012).
2. Stergachis, A. B. *et al.* XDevelopmental fate and cellular maturity encoded in human regulatory DNA landscapes. *Cell* **154**, (2013).
3. Brunner, A. *et al.* Ultra-high sensitivity mass spectrometry quantifies single-cell proteome changes upon perturbation. *Mol Syst Biol* **18**, (2022).
4. Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell. *Nat Methods* **6**, (2009).
5. Lareau, C. A. *et al.* Massively parallel single-cell mitochondrial DNA genotyping and chromatin profiling. *Nat Biotechnol* **39**, (2021).
6. Zenobi, R. Single-cell metabolomics: Analytical and biological perspectives. *Science* vol. 342 Preprint at <https://doi.org/10.1126/science.1243259> (2013).
7. Boudreau, N. & Bissell, M. J. Extracellular matrix signaling: Integration of form and function in normal and malignant cells. *Curr Opin Cell Biol* **10**, (1998).
8. Discher, D. E., Janmey, P. & Wang, Y. L. Tissue cells feel and respond to the stiffness of their substrate. *Science* vol. 310 Preprint at <https://doi.org/10.1126/science.1116995> (2005).
9. Heldin, C. H., Lu, B., Evans, R. & Gutkind, J. S. Signals and receptors. *Cold Spring Harb Perspect Biol* **8**, (2016).
10. Colombo, M., Raposo, G. & Théry, C. Biogenesis, secretion, and intercellular interactions of exosomes and other extracellular vesicles. *Annual review of cell and developmental biology* vol. 30 Preprint at <https://doi.org/10.1146/annurev-cellbio-101512-122326> (2014).
11. McFaline-Figueroa, J. L. *et al.* A pooled single-cell genetic screen identifies regulatory checkpoints in the continuum of the epithelial-to-mesenchymal transition. *Nature Genetics* **2019 51:9 51**, 1389–1398 (2019).
12. Deshmukh, A. P. *et al.* Identification of EMT signaling cross-talk and gene regulatory networks by single-cell RNA sequencing. *Proceedings of the National Academy of Sciences* **118**, e2102050118 (2021).
13. Vandereyken, K., Sifrim, A., Thienpont, B. & Thierry Voet, &. Methods and applications for single-cell and spatial multi-omics. *Nature Reviews Genetics* | **24**, 494–515 (2023).
14. Giesen, C. *et al.* Highly multiplexed imaging of tumor tissues with subcellular resolution by mass cytometry. *Nature Methods* **2014 11:4 11**, 417–422 (2014).
15. Lin, J. R., Fallahi-Sichani, M. & Sorger, P. K. Highly multiplexed imaging of single cells using a high-throughput cyclic immunofluorescence method. *Nature Communications* **2015 6:1 6**, 1–7 (2015).
16. Tsujikawa, T. *et al.* Quantitative Multiplex Immunohistochemistry Reveals Myeloid-Inflamed Tumor-Immune Complexity Associated with Poor Prognosis. *Cell Rep* **19**, 203–217 (2017).
17. C. Greenwald, E., Mehta, S. & Zhang, J. Genetically Encoded Fluorescent Biosensors Illuminate the Spatiotemporal Regulation of Signaling Networks. *Chem Rev* **118**, 11707–11794 (2018).
18. Linghu, C. *et al.* Spatial Multiplexing of Fluorescent Reporters for Imaging Signaling Network Dynamics. *Cell* **183**, 1682-1698.e24 (2020).

19. Davies, A. E. *et al.* Article Systems-Level Properties of EGFR-RAS-ERK Signaling Amplify Local Signals to Generate Dynamic Gene Expression Heterogeneity. *Cell Syst* **11**, 161-175.e5 (2020).
20. Yang, J. M. *et al.* Deciphering cell signaling networks with massively multiplexed biosensor barcoding. *Cell* **184**, 6193-6206.e14 (2021).
21. Neumann, B. *et al.* High-throughput RNAi screening by time-lapse imaging of live human cells. *Nat Methods* **3**, 385–390 (2006).
22. Held, M. *et al.* CellCognition: Time-resolved phenotype annotation in high-throughput live cell imaging. *Nat Methods* **7**, 747–754 (2010).
23. Nketia, T. A., Sailem, H., Rohde, G., Machiraju, R. & Rittscher, J. Analysis of live cell images: Methods, tools and opportunities. *Methods* vol. 115 65–79 Preprint at <https://doi.org/10.1016/j.ymeth.2017.02.007> (2017).
24. Stirling, D. R. *et al.* CellProfiler 4: improvements in speed, utility and usability. *BMC Bioinformatics* **22**, 1–11 (2021).
25. Copperman, J., Gross, S. M., Chang, Y. H., Heiser, L. M. & Zuckerman, D. M. Morphodynamical cell state description via live-cell imaging trajectory embedding. *Commun Biol* **6**, (2023).
26. Karacosta, L. G. *et al.* Mapping lung cancer epithelial-mesenchymal transition states and trajectories with single-cell resolution. *Nat Commun* **10**, 1–15 (2019).
27. Cho, H., Kuo, Y. H. & Rockne, R. C. Comparison of cell state models derived from single-cell RNA sequencing data: graph versus multi-dimensional space. *Mathematical Biosciences and Engineering* **19**, (2022).
28. Tong, A. *et al.* Learning transcriptional and regulatory dynamics driving cancer cell plasticity using neural ODE-based optimal transport. *bioRxiv* (2023).
29. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, (2014).
30. Bendall, S. C. *et al.* Single-cell trajectory detection uncovers progression and regulatory coordination in human b cell development. *Cell* **157**, 714–725 (2014).
31. La Manno, G. *et al.* RNA velocity of single cells. *Nature* **560**, 494–498 (2018).
32. Bergen, V., Lange, M., Peidli, S., Wolf, F. A. & Theis, F. J. Generalizing RNA velocity to transient cell states through dynamical modeling. *Nat Biotechnol* **38**, 1408–1414 (2020).
33. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods* **12**, (2015).
34. Wang, Z. *et al.* Transcriptome Deconvolution of Heterogeneous Tumor Samples with Immune Infiltration. *iScience* **9**, (2018).
35. Wang, X., Park, J., Susztak, K., Zhang, N. R. & Li, M. Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. *Nat Commun* **10**, (2019).
36. Danziger, S. A. *et al.* AdApTS: Automated deconvolution augmentation of profiles for tissue specific cells. *PLoS One* **14**, (2019).
37. Noé, F. & Fischer, S. Transition networks for modeling the kinetics of conformational change in macromolecules. *Curr Opin Struct Biol* **18**, 154–162 (2008).
38. Husic, B. E. & Pande, V. S. Markov State Models: From an Art to a Science. *Journal of the American Chemical Society* vol. 140 Preprint at <https://doi.org/10.1021/jacs.7b12191> (2018).
39. Wang, W. *et al.* Live-cell imaging and analysis reveal cell phenotypic transition dynamics inherently missing in snapshot data. *Sci Adv* **6**, eaba9319 (2020).

40. Wang, W., Poe, D., Yang, Y., Hyatt, T. & Xing, J. Epithelial-to-mesenchymal transition proceeds through directional destabilization of multidimensional attractor. *Elife* **11**, (2022).
41. Soule, H. D. *et al.* Isolation and Characterization of a Spontaneously Immortalized Human Breast Epithelial Cell Line, MCF-10. *Cancer Res* **50**, (1990).
42. Witt, A. E. *et al.* Functional proteomics approach to investigate the biological activities of cDNAs implicated in breast cancer. *J Proteome Res* **5**, (2006).
43. Melani, M., Simpson, K. J., Brugge, J. S. & Montell, D. Regulation of Cell Adhesion and Collective Cell Migration by Hindsight and Its Human Homolog RREB1. *Current Biology* **18**, (2008).
44. Seton-Rogers, S. E. *et al.* Cooperation of the ErbB2 receptor and transforming growth factor  $\beta$  in induction of migration and invasion in mammary epithelial cells. *Proc Natl Acad Sci U S A* **101**, (2004).
45. Debnath, J. *et al.* The role of apoptosis in creating and maintaining luminal space within normal and oncogene-expressing mammary acini. *Cell* **111**, (2002).
46. Debnath, J., Muthuswamy, S. K. & Brugge, J. S. Morphogenesis and oncogenesis of MCF-10A mammary epithelial acini grown in three-dimensional basement membrane cultures. *Methods* vol. 30 Preprint at [https://doi.org/10.1016/S1046-2023\(03\)00032-X](https://doi.org/10.1016/S1046-2023(03)00032-X) (2003).
47. Sampattavanich, S. *et al.* Encoding Growth Factor Identity in the Temporal Dynamics of FOXO3 under the Combinatorial Control of ERK and AKT Kinases. *Cell Syst* **6**, (2018).
48. Caldera, M. *et al.* Mapping the perturbome network of cellular perturbations. *Nat Commun* **10**, (2019).
49. Gross, S. M. *et al.* A multi-omic analysis of MCF10A cells provides a resource for integrative assessment of ligand-mediated molecular and phenotypic responses. *Commun Biol* **5**, (2022).
50. Espinosa-Neira, R., Mejia-Rangel, J., Cortes-Reynosa, P. & Salazar, E. P. Linoleic acid induces an EMT-like process in mammary epithelial cells MCF10A. *International Journal of Biochemistry and Cell Biology* **43**, (2011).
51. Galindo-Hernandez, O., Serna-Marquez, N., Castillo-Sanchez, R. & Salazar, E. P. Extracellular vesicles from MDA-MB-231 breast cancer cells stimulated with linoleic acid promote an EMT-like process in MCF10A cells. *Prostaglandins Leukot Essent Fatty Acids* **91**, (2014).
52. Zhang, J. *et al.* TGF- $\beta$ -induced epithelial-to-mesenchymal transition proceeds through stepwise activation of multiple feedback loops. *Sci Signal* **7**, (2014).
53. Mori, S. *et al.* Enhanced expression of integrin  $\alpha\beta 3$  induced by TGF- $\beta$  is required for the enhancing effect of fibroblast growth factor 1 (FGF1) in TGF- $\beta$ -induced epithelial-mesenchymal transition (EMT) in mammary epithelial cells. *PLoS One* **10**, (2015).
54. Rodriguez-Monteros, C., Díaz-Aragon, R., Leal-Orta, E., Cortes-Reynosa, P. & Perez Salazar, E. Insulin induces an EMT-like process in mammary epithelial cells MCF10A. *J Cell Biochem* **119**, (2018).
55. Olea-Flores, M. *et al.* Leptin promotes expression of EMT-related transcription factors and invasion in a src and FAK-dependent pathway in MCF10a mammary epithelial cells. *Cells* **8**, (2019).



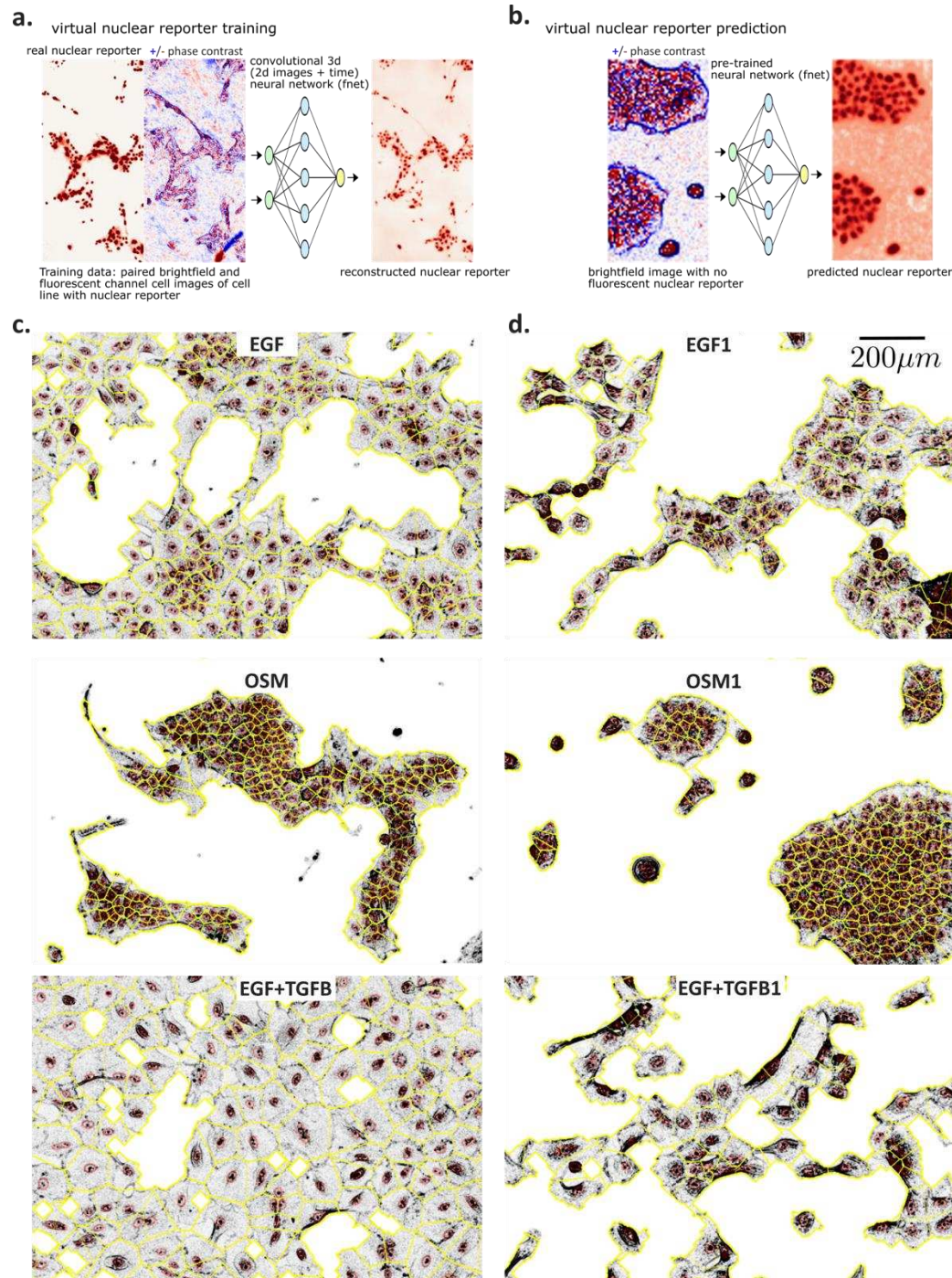
56. Antón-García, P. *et al.* TGFβ1-Induced EMT in the MCF10A Mammary Epithelial Cell Line Model Is Executed Independently of SNAIL1 and ZEB1 but Relies on JUNB-Coordinated Transcriptional Regulation. *Cancers (Basel)* **15**, (2023).
57. Paul, I. *et al.* Parallelized multidimensional analytic framework applied to mammary epithelial cells uncovers regulatory principles in EMT. *Nat Commun* **14**, (2023).
58. Spencer, S. L. *et al.* XThe proliferation-quiescence decision is controlled by a bifurcation in CDK2 activity at mitotic exit. *Cell* **155**, (2013).
59. E, W. & Vanden-Eijnden, E. Transition-path theory and path-finding algorithms for the study of rare events. *Annu Rev Phys Chem* **61**, (2010).
60. Scherer, M. K. *et al.* PyEMMA 2: A Software Package for Estimation, Validation, and Analysis of Markov Models. *J Chem Theory Comput* **11**, (2015).
61. Ounkomol, C., Seshamani, S., Maleckar, M. M., Collman, F. & Johnson, G. R. Label-free prediction of three-dimensional fluorescence images from transmitted-light microscopy. *Nat Methods* **15**, 917–920 (2018).
62. Ulicna, K., Vallardi, G., Charras, G. & Lowe, A. R. Automated Deep Lineage Tree Analysis Using a Bayesian Single Cell Tracking Approach. *Front Comput Sci* **3**, (2021).
63. Röblitz, S. & Weber, M. Fuzzy spectral clustering by PCCA+: Application to Markov state models and data classification. *Adv Data Anal Classif* **7**, (2013).
64. Reuter, B., Weber, M., Fackeldey, K., Röblitz, S. & Garcia, M. E. Generalized Markov State Modeling Method for Nonequilibrium Biomolecular Dynamics: Exemplified on Amyloid β Conformational Dynamics Driven by an Oscillating Electric Field. *J Chem Theory Comput* **14**, 3579–3594 (2018).
65. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *ArXiv* (2018).
66. Dongre, A. & Weinberg, R. A. New insights into the mechanisms of epithelial–mesenchymal transition and implications for cancer. *Nature Reviews Molecular Cell Biology* vol. 20 Preprint at <https://doi.org/10.1038/s41580-018-0080-4> (2019).
67. Donovan, J. & Slingerland, J. Transforming growth factor-β and breast cancer: Cell cycle arrest by transforming growth factor-β and its disruption in cancer. *Breast Cancer Research* vol. 2 Preprint at <https://doi.org/10.1186/bcr43> (2000).
68. Mejlvang, J. *et al.* Direct repression of cyclin D1 by SIP1 attenuates cell cycle progression in cells undergoing an epithelial mesenchymal transition. *Mol Biol Cell* **18**, (2007).
69. Lovisa, S. *et al.* Epithelial-to-mesenchymal transition induces cell cycle arrest and parenchymal damage in renal fibrosis. *Nat Med* **21**, (2015).
70. Kohrman, A. Q. & Matus, D. Q. Divide or Conquer: Cell Cycle Regulation of Invasive Behavior. *Trends in Cell Biology* vol. 27 Preprint at <https://doi.org/10.1016/j.tcb.2016.08.003> (2017).
71. Petersen, O. W., Rønnov-Jessen, L., Howlett, A. R. & Bissell, M. J. Interaction with basement membrane serves to rapidly distinguish growth and differentiation pattern of normal and malignant human breast epithelial cells. *Proceedings of the National Academy of Sciences* **89**, 9064–9068 (1992).
72. Battle, C. *et al.* Broken detailed balance at mesoscopic scales in active biological systems. *Science (1979)* **352**, 604–607 (2016).
73. Kimmel, J. C., Chang, A. Y., Brack, A. S. & Marshall, W. F. Inferring cell state by quantitative motility analysis reveals a dynamic state system and broken detailed balance. *PLoS Comput Biol* **14**, e1005927 (2018).

74. Liberzon, A. *et al.* The Molecular Signatures Database Hallmark Gene Set Collection. *Cell Syst* **1**, (2015).
75. Zaritsky, A. *et al.* Interpretable deep learning uncovers cellular properties in label-free live cell images that are predictive of highly metastatic melanoma. *Cell Syst* **12**, 733-747.e6 (2021).
76. Nguyen, P. *et al.* Unsupervised discovery of dynamic cell phenotypic states from transmitted light movies. *PLoS Comput Biol* **17**, (2021).
77. Chow, Y. L., Singh, S., Carpenter, A. E. & Way, G. P. Predicting drug polypharmacology from cell morphology readouts using variational autoencoder latent space arithmetic. *PLoS Comput Biol* **18**, (2022).
78. Ternes, L. *et al.* A multi-encoder variational autoencoder controls multiple transformational features in single-cell image analysis. *Communications Biology* **2022 5:1** **5**, 1–10 (2022).
79. Schau, G., Burlingame, E. & Chang, Y. H. DISSECT: DISentangle Sharable ConTent for Multimodal Integration and Crosswise-mapping. in *Proceedings of the IEEE Conference on Decision and Control* vols 2020-December 5092–5097 (Institute of Electrical and Electronics Engineers Inc., 2020).
80. Pratt, L. R. A statistical method for identifying transition states in high dimensional problems. *J Chem Phys* **85**, (1986).
81. Bolhuis, P. G., Chandler, D., Dellago, C. & Geissler, P. L. Transition Path Sampling: Throwing ropes over rough mountain passes, in the dark. *Annu Rev Phys Chem* **53**, 291–318 (2002).
82. Chong, L. T., Saglam, A. S. & Zuckerman, D. M. Path-sampling strategies for simulating rare events in biomolecular systems. *Current Opinion in Structural Biology* vol. 43 88–94 Preprint at <https://doi.org/10.1016/j.sbi.2016.11.019> (2017).
83. Yang, J. & Weinberg, R. A. Epithelial-Mesenchymal Transition: At the Crossroads of Development and Tumor Metastasis. *Developmental Cell* vol. 14 Preprint at <https://doi.org/10.1016/j.devcel.2008.05.009> (2008).
84. Grasset, E. M. *et al.* Triple-negative breast cancer metastasis involves complex epithelial-mesenchymal transition dynamics and requires vimentin. *Sci Transl Med* **14**, (2022).
85. Risom, T. *et al.* Differentiation-state plasticity is a targetable resistance mechanism in basal-like breast cancer. *Nat Commun* **9**, 3815 (2018).
86. Gambardella, G. *et al.* A single-cell analysis of breast cancer cell lines to study tumour heterogeneity and drug response. *Nat Commun* **13**, (2022).
87. Tatarova, Z. *et al.* A multiplex implantable microdevice assay identifies synergistic combinations of cancer immunotherapies and conventional drugs. *Nat Biotechnol* **40**, (2022).
88. Pong, A., Mah, C. K., Yeo, G. W. & Lewis, N. E. Computational cell–cell interaction technologies drive mechanistic and biomarker discovery in the tumor microenvironment. *Curr Opin Biotechnol* (2024).
89. Sauro, H. M. & Kholodenko, B. N. Quantitative analysis of signaling networks. *Prog Biophys Mol Biol* **86**, (2004).
90. Alon, U. Network motifs: Theory and experimental approaches. *Nature Reviews Genetics* vol. 8 Preprint at <https://doi.org/10.1038/nrg2102> (2007).



91. Aldridge, B. B., Burke, J. M., Lauffenburger, D. A. & Sorger, P. K. Physicochemical modelling of cell signalling pathways. *Nature Cell Biology* vol. 8 Preprint at <https://doi.org/10.1038/ncb1497> (2006).
92. Birtwistle, M. R. *et al.* Ligand-dependent responses of the ErbB signaling network: Experimental and modeling analyses. *Mol Syst Biol* **3**, (2007).
93. Gross, S. M. *et al.* Analysis and modeling of cancer drug responses using cell cycle phase-specific rate effects. *Nat Commun* **14**, (2023).
94. Sommer, C., Straehle, C., Kothe, U. & Hamprecht, F. A. Ilastik: Interactive learning and segmentation toolkit. in *Proceedings - International Symposium on Biomedical Imaging* 230–233 (2011). doi:10.1109/ISBI.2011.5872394.
95. Thévenaz, P., Ruttimann, U. E. & Unser, M. A pyramid approach to subpixel registration based on intensity. *IEEE Transactions on Image Processing* **7**, 27–41 (1998).
96. Coelho, L. P. Mahotas: Open source software for scriptable computer vision. *J Open Res Softw* **1**, e3 (2013).
97. Alizadeh, E., Xu, W., Castle, J., Foss, J. & Prasad, A. TISMorph: A tool to quantify texture, irregularity and spreading of single cells. *PLoS One* **14**, e0217346 (2019).
98. Bertrand, T. *et al.* Clustering and ordering in cell assemblies with generic asymmetric aligning interactions. *ArXiv* (2020).
99. Atev, S. E. Using Asymmetry in the Spectral Clustering of Trajectories. (University of Minnesota, 2011).
100. Nüske, F., Keller, B. G., Pérez-Hernández, G., Mey, A. S. J. S. & Noé, F. Variational approach to molecular kinetics. *J Chem Theory Comput* **10**, 1739–1752 (2014).
101. Fang, Z., Liu, X. & Peltz, G. GSEAPy: a comprehensive package for performing gene set enrichment analysis in Python. *Bioinformatics* **39**, (2023).

# Supplementary Data Tables and Figures



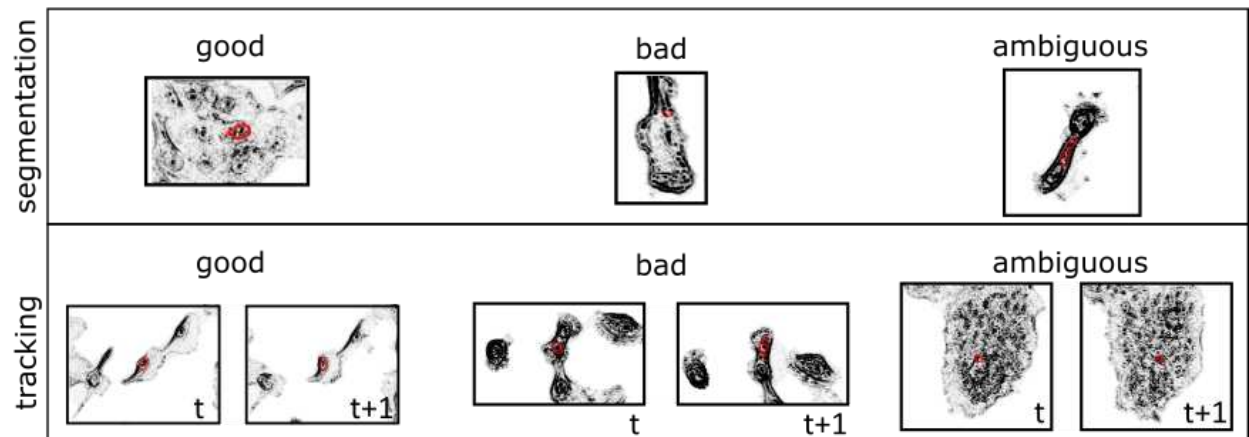
## Supplementary Figure 1: Virtual nuclear staining and nuclear center-based segmentation

a) Paired nuclear reporter (red) and z-normalized phase contrast (+ red - blue) training data input, and reconstructed nuclear reporter output b) Nuclear reporter prediction from out of sample phase contrast images (OSM condition). c) Examples of nuclear segmentations (red) and associated Voronoi boundaries (yellow) overlaid upon phase contrast images (z-normalized absolute value, gray)

**Supplementary Data Table 1. Segmentation and tracking manual validation.** 100 cells per treatment were randomly selected, and evaluated by eye to qualitatively assess segmentation and tracking accuracy. Fraction segmented was estimated by the image area covered by segmented masks divided by the area selected as being occupied by cells from the ilastik random forest pixel classifier. Segmentation performance from dataset 1 (e.g. EGF1, HGF1) is decreased because for these data the cells did not express a nuclear reporter, and the nucleus was detected via the virtual staining approach only. Tracking performance is decreased as well, due to the decreased segmentation performance and increased time between frames (30 minutes as compared to 15 minutes).

ligand	PBS	EGF	OSM	TGFB	TGFB +EGF	OSM +EGF	TGFB +OSM	TGFB +OSM +EGF
count/manual	99%	99%	98%	96%	99%	98%	96%	99%
% good seg	92%	97%	97%	95%	96%	96%	95%	95%
% bad seg	1%	1%	1%	1%	2%	2%	0%	1%
%ambiguous seg	7%	2%	2%	4%	2%	2%	5%	4%
% tracked	97%	91%	93%	92%	95%	93%	88%	90%
% good tracks	100%	100%	98%	100%	99%	99%	100%	100%
% bad tracks	0%	0%	2%	0%	0%	0%	0%	0%
%ambiguous tracks	0%	0%	0%	0%	1%	1%	0%	0%

ligand	PBS1	EGF1	HGF1	OSM1	IFNG+ EGF1	BMP2+ EGF1	EGF+ TGFB1
counts/manual	79%	94%	74%	90%	99%	102%	113%
% good seg	80%	62%	74%	86%	68%	78%	75%
% bad seg	11%	25%	11%	4%	24%	16%	9%
%ambiguous seg	9%	10%	14%	10%	8%	6%	16%
% tracked	78%	52%	78%	69%	55%	66%	49%
% good tracks	94%	94%	99%	88%	91%	95%	96%
% bad tracks	4%	6%	1%	7%	5%	5%	4%
%ambiguous tracks	2%	0%	0%	5%	4%	0%	0%



### Supplementary Figure 2. Segmentation and tracking manual validation examples

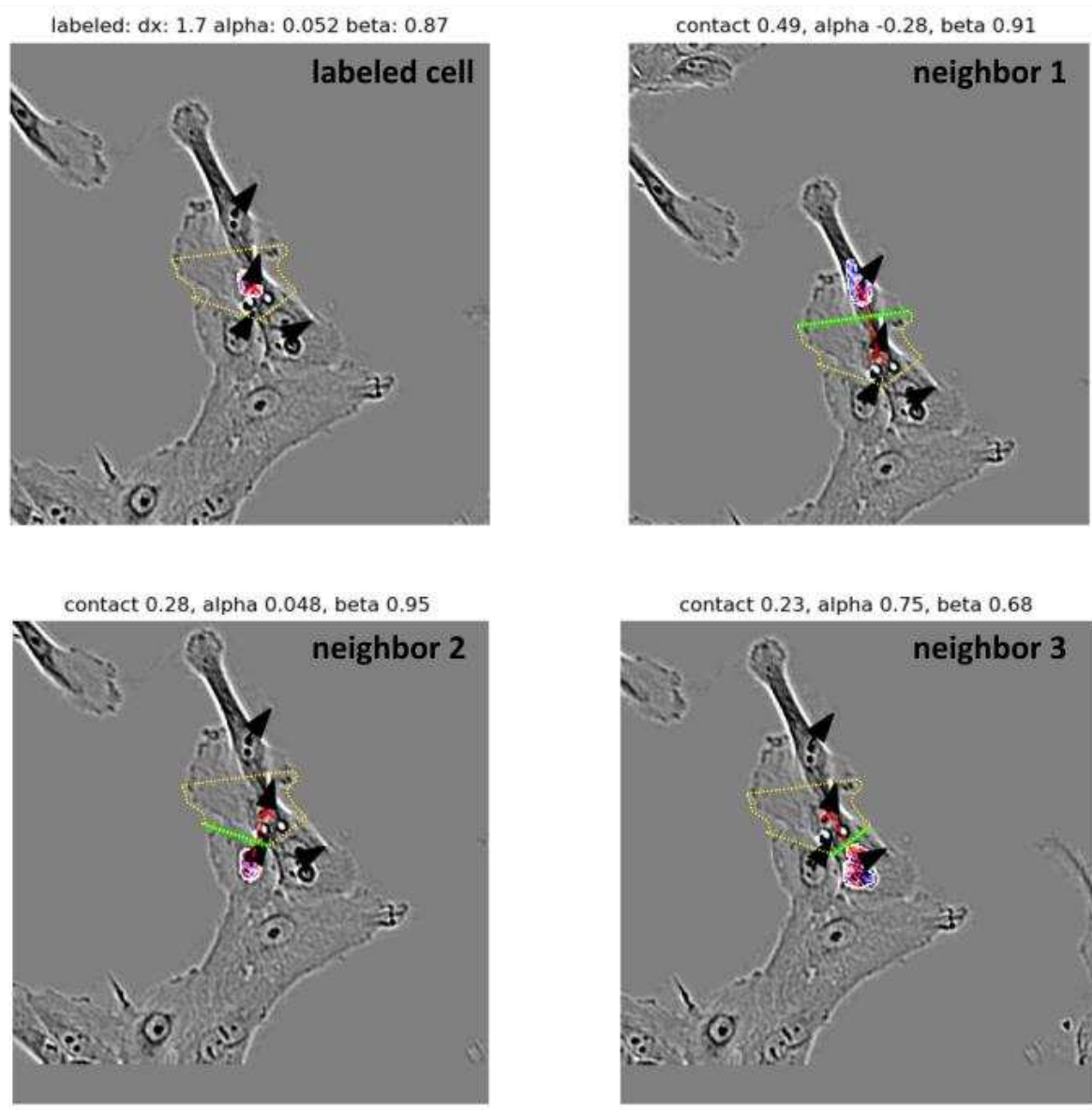
Examples of good, bad, and ambiguous qualitative validation categories for segmentation and tracking. Absolute value of z-normalized phase contrast in gray, nuclear segmentation in red.

**Supplementary Data Table 2. Test set gene expression validation with trajectory length.**

test treatment	traj length (hrs)	pred. to exp. corr	diff. from null $\rho - \rho_{\text{null}}$	upreg true pos. rate	downreg true pos. rate
OSM+EGF	0	.26	-0.38	.51	.80
EGF+TGFB	0	.76	0.09	.63	.89
OSM+EGF+TGFB	0	.77	0.04	.77	.95
OSM+EGF	1	.51	0.02	.55	<b>.99</b>
EGF+TGFB	1	.76	0.22	.68	.93
OSM+EGF+TGFB	1	.75	0.14	.77	.93
OSM+EGF	4	.50	-0.22	.51	.98
EGF+TGFB	4	.74	0.05	.68	.80
OSM+EGF+TGFB	4	.73	-0.04	.77	<b>.96</b>
OSM+EGF	8	.55	-0.15	.52	<b>.99</b>
EGF+TGFB	8	.69	0	.68	<b>.87</b>
OSM+EGF+TGFB	8	.66	-0.08	.77	.94
OSM+EGF	10	<b>.70</b>	<b>0.21</b>	<b>.72</b>	.81
EGF+TGFB	10	<b>.81</b>	<b>0.27</b>	<b>.72</b>	.84
OSM+EGF+TGFB	10	.77	<b>0.12</b>	<b>.85</b>	.81
OSM+EGF	12	.71	0.2	.70	.77
EGF+TGFB	12	.74	0.14	.69	.85
OSM+EGF+TGFB	12	.64	0.01	.80	.80
OSM+EGF	16	.66	-0.03	.65	.92
EGF+TGFB	16	.79	0.15	.70	.85
OSM+EGF+TGFB	16	<b>.79</b>	0.09	.80	.88



$$f_{\text{labeled}} = .49f_1 + .28f_2 + .23f_3$$

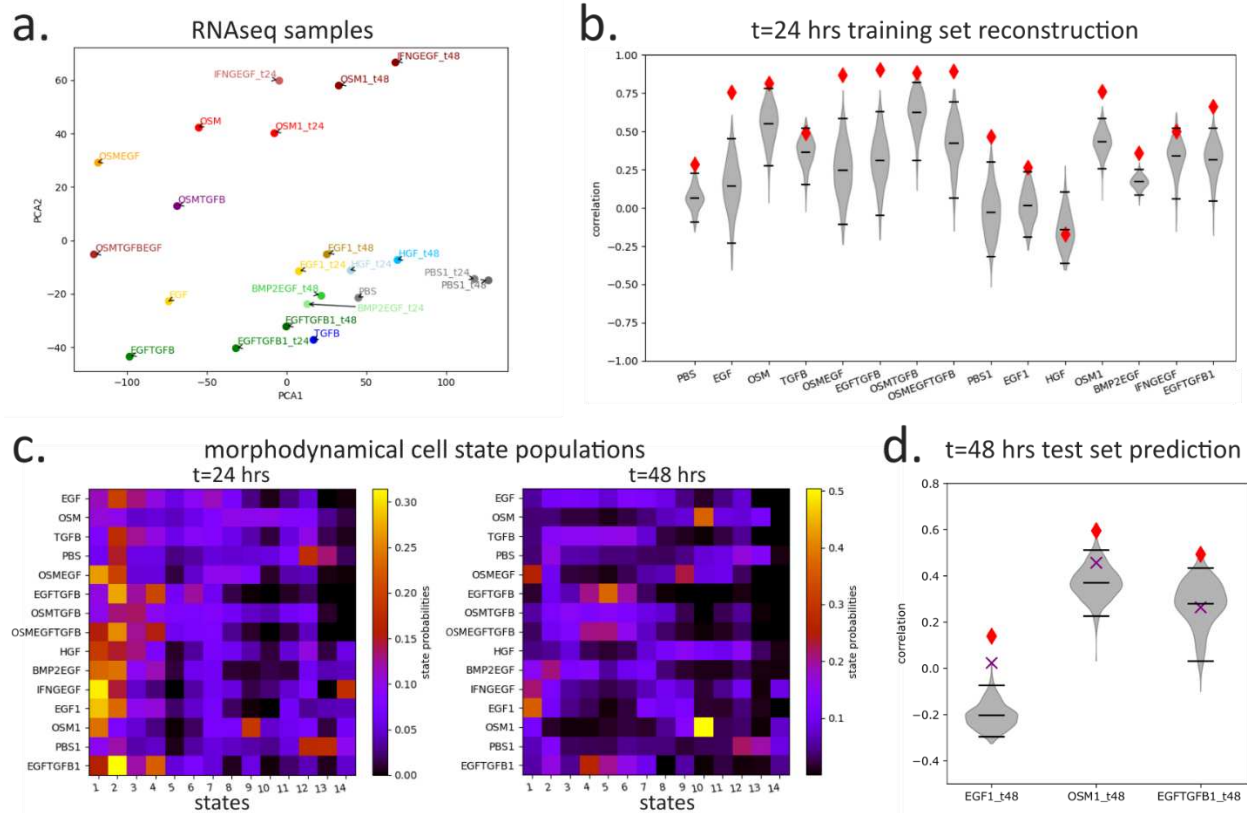


$\rightarrow \Delta \vec{x}$  (dt = 15min)  
 labeled cell boundary  
 neighbor to labeled separation vector  
 shared cell boundary

### Supplementary Figure 3: Single-cell and neighborhood motility feature

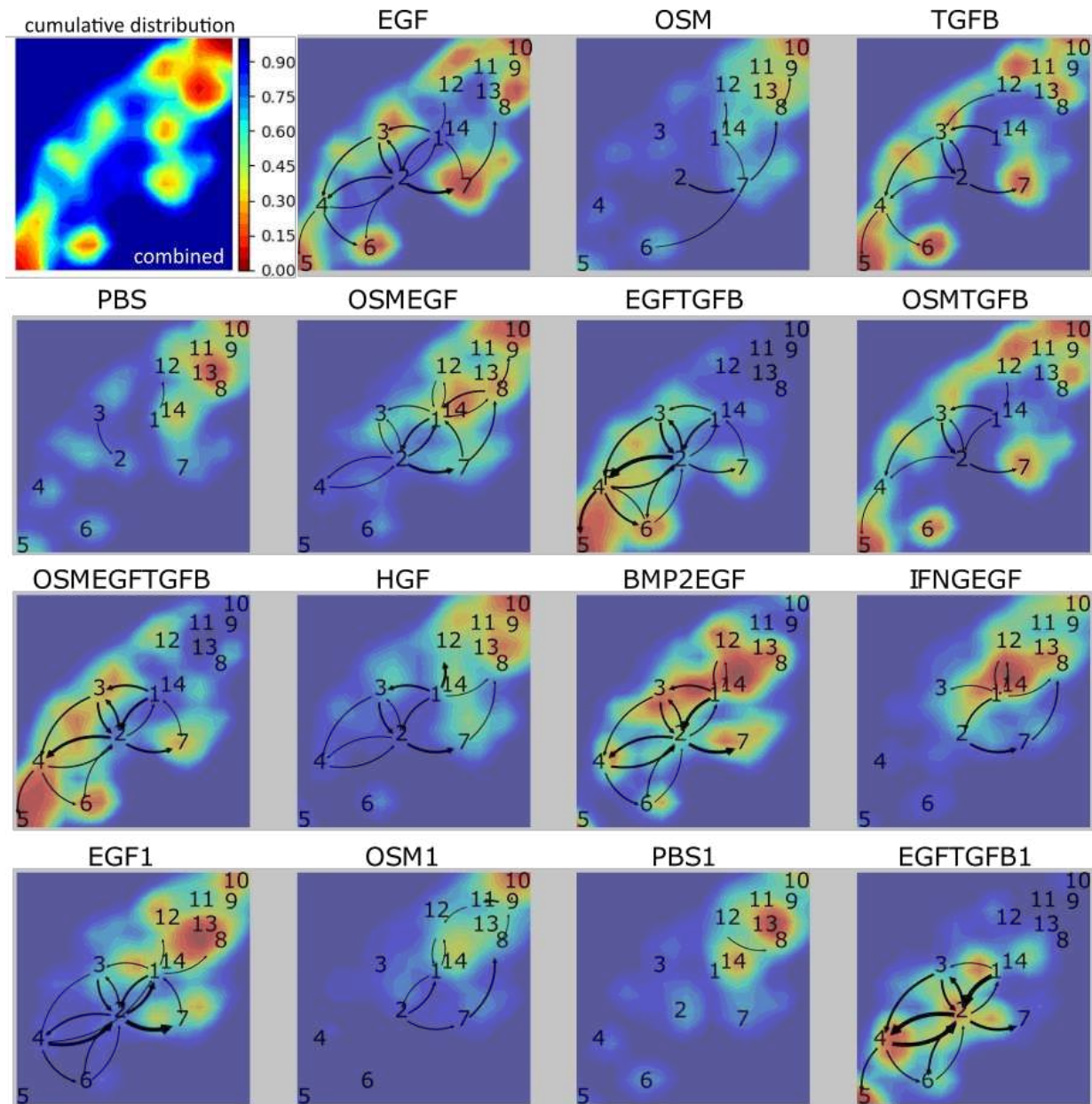
Single-cell featurization of motility for the labeled cell in the upper right (nuclei in color and Voronoi segmentation in yellow) taken as the magnitude of the displacement from previous frame. The single-cell motility in the context of its local neighborhood taken as the neighbor-weighted average of the 3 boundary cells (upper right, lower left, lower right).





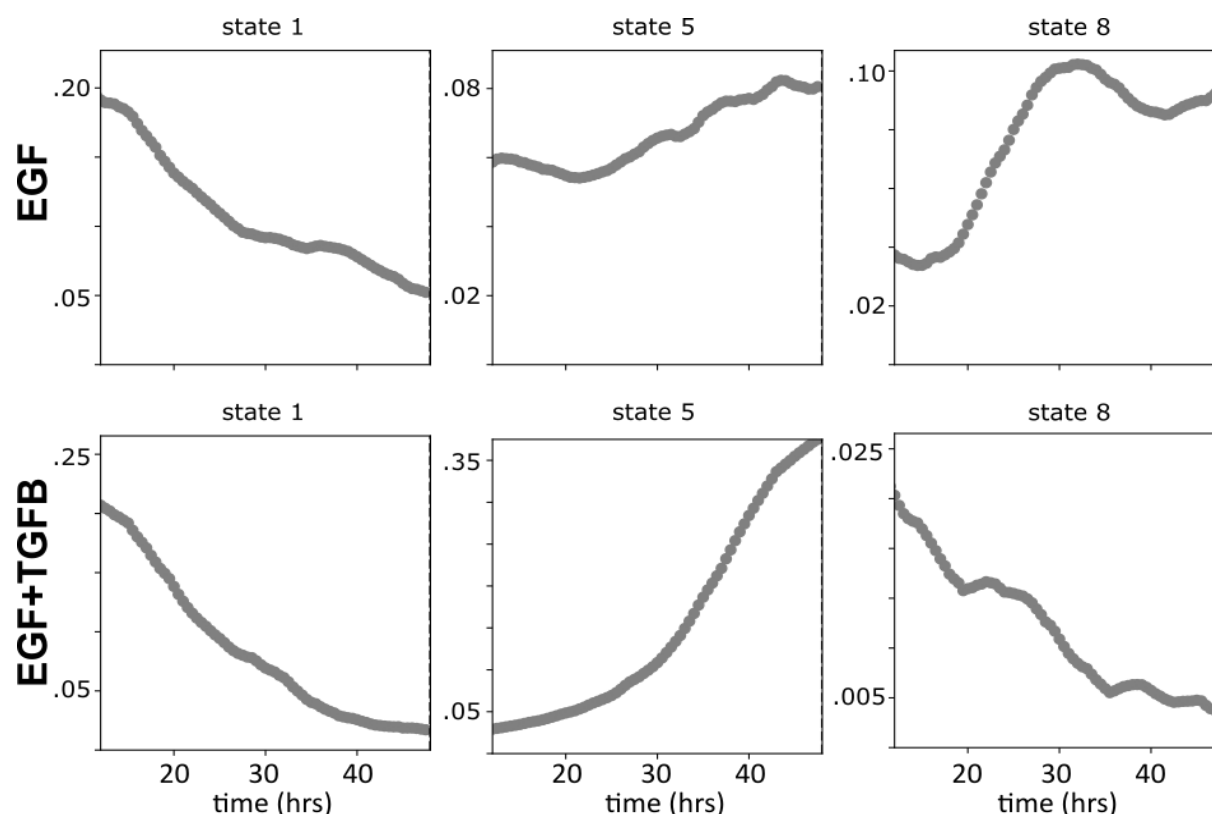
### Supplementary Figure 4: Bulk RNAseq decomposition and time dependence

a) PCA1/2 projection of the RNAseq differential expression, showing sorting by ligand treatment and timepoint. b) Correlation between training set reconstruction and real experimental differential expression. c) Morphodynamical state populations at t=24, 48 hrs d) Test set prediction of RNAseq at 48 hours using t=24 hrs trained morphodynamical state gene expression profiles and measured live-cell state populations at t=48hrs.



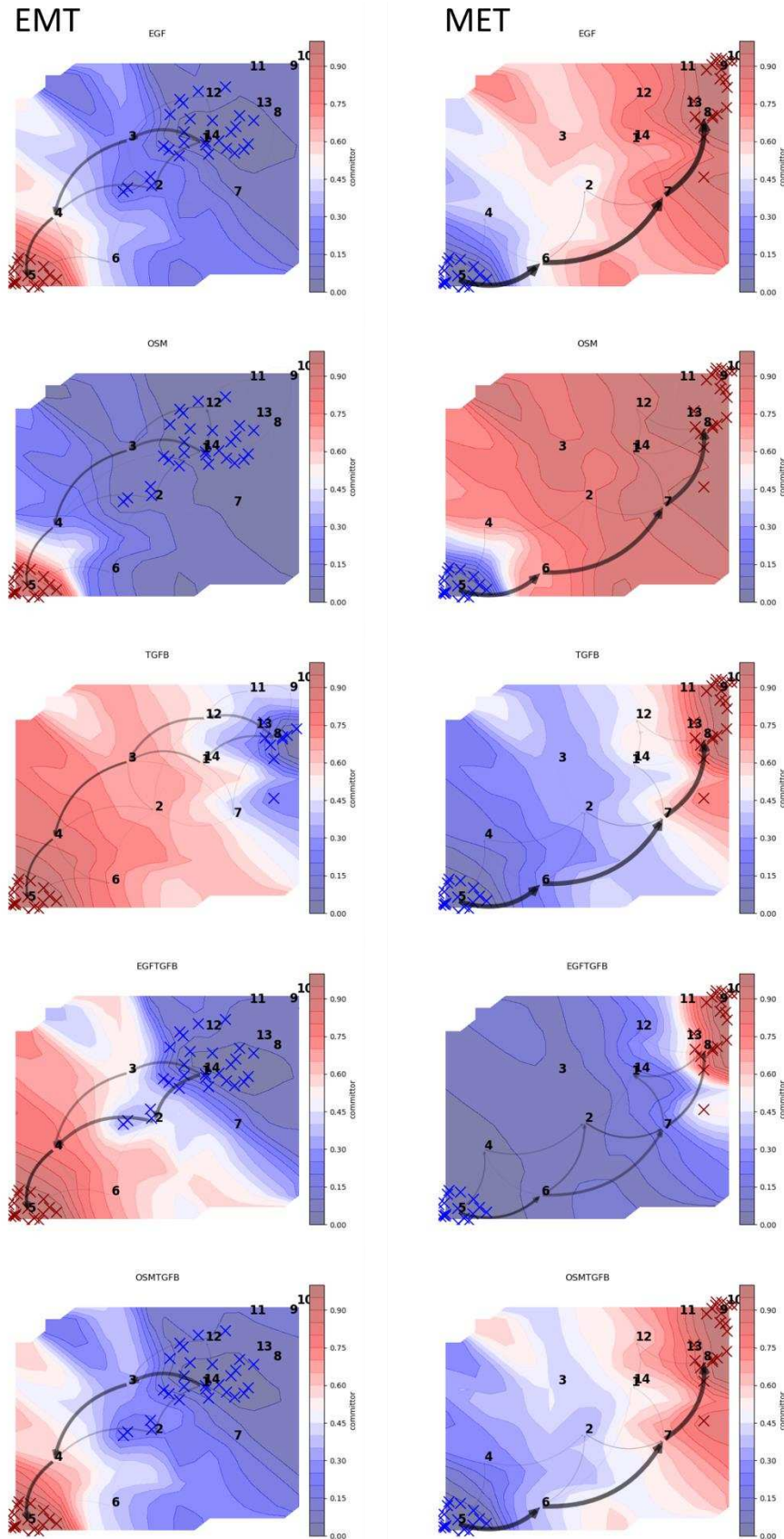
### Supplementary Figure 5: Ligand-dependent populations and cell state flows

Cumulative populations in the UMAP embedding space (blue to red), and state-state transition flows at t=24hrs, in each ligand treatment.



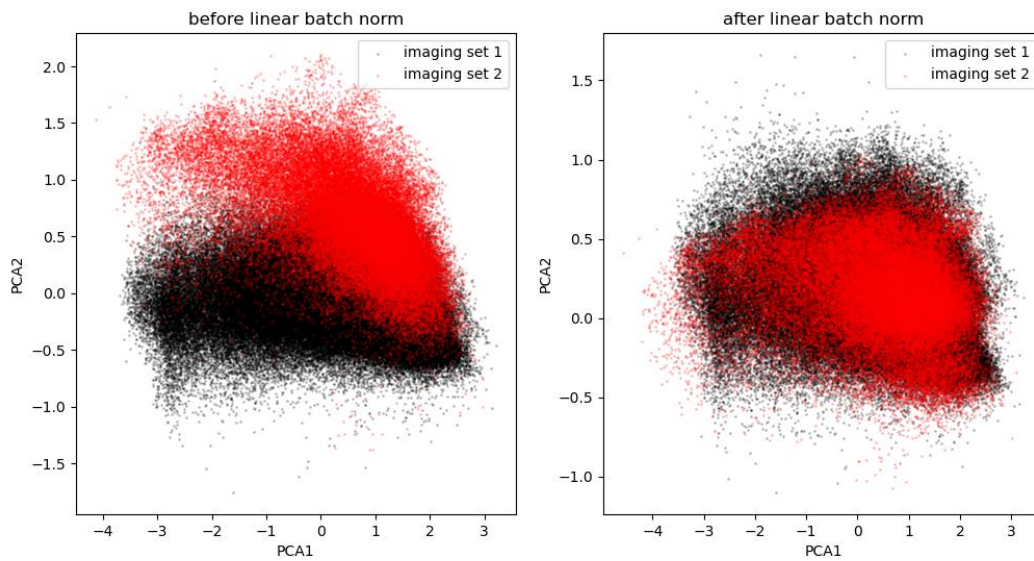
### Supplementary Figure 6: EMT/MET initial and final state probabilities

Live-cell imaging inferred initial and final EMT/MET state populations as a function of time. EMT initial state 1 depopulates over time (rightmost plots), while EMT final state 5 increases over time in EGF+TGFB (lower middle), while MET final state 8 increases over time in EGF conditions (upper right).



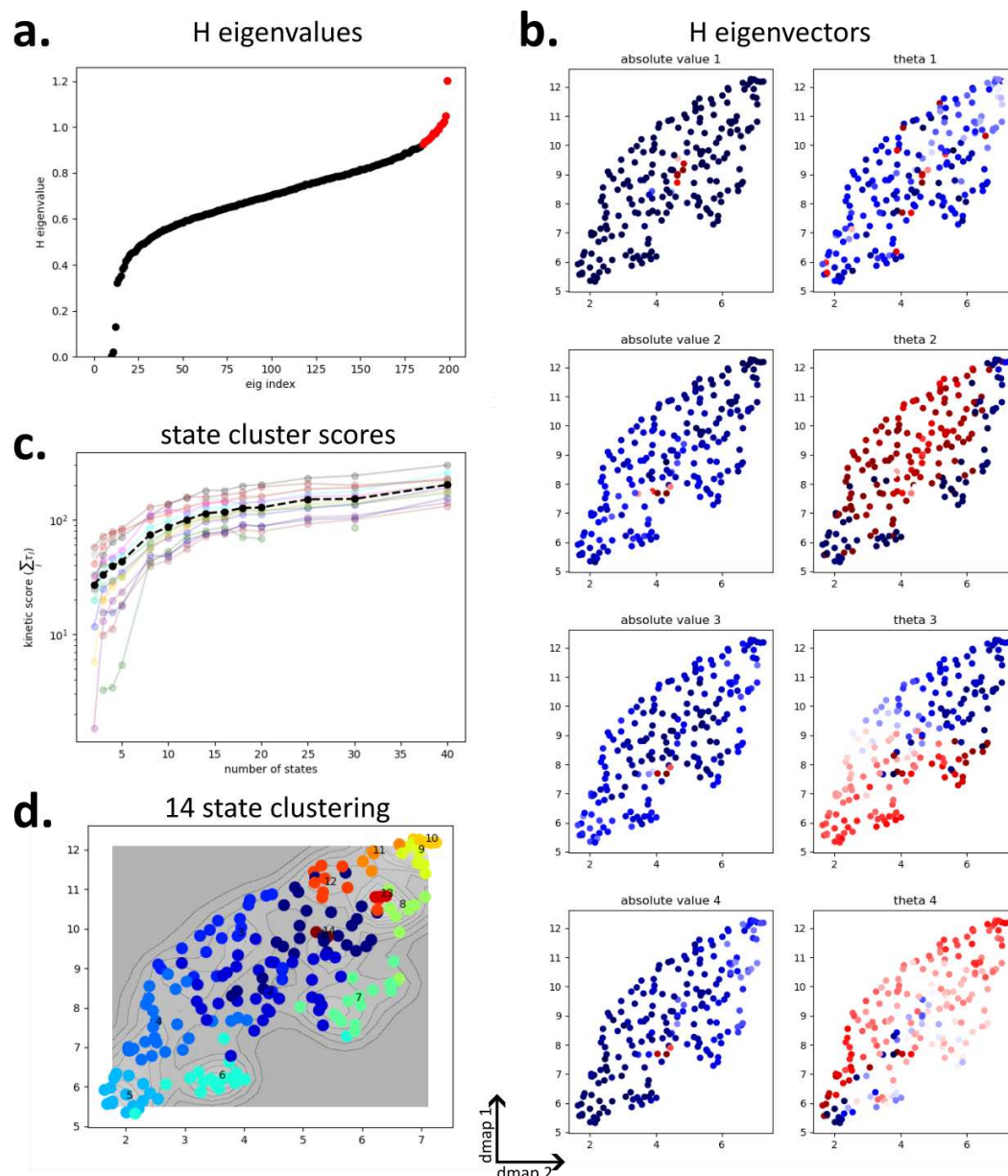
**Supplementary Figure 7: EMT/MET morphodynamical cell state change sequences by ligand treatment**  
Possible EMT cell state change sequences (initial state 1, final state 5) left, and MET cell state change sequences (initial state 5, final state 8) on the right (black arrows, thickness proportional to transition flux), with final state commitment probability (blue to red) calculated from the 200 k-means state centers and averaged over the UMAP surface.





### Supplementary Figure 8: Feature batch normalization

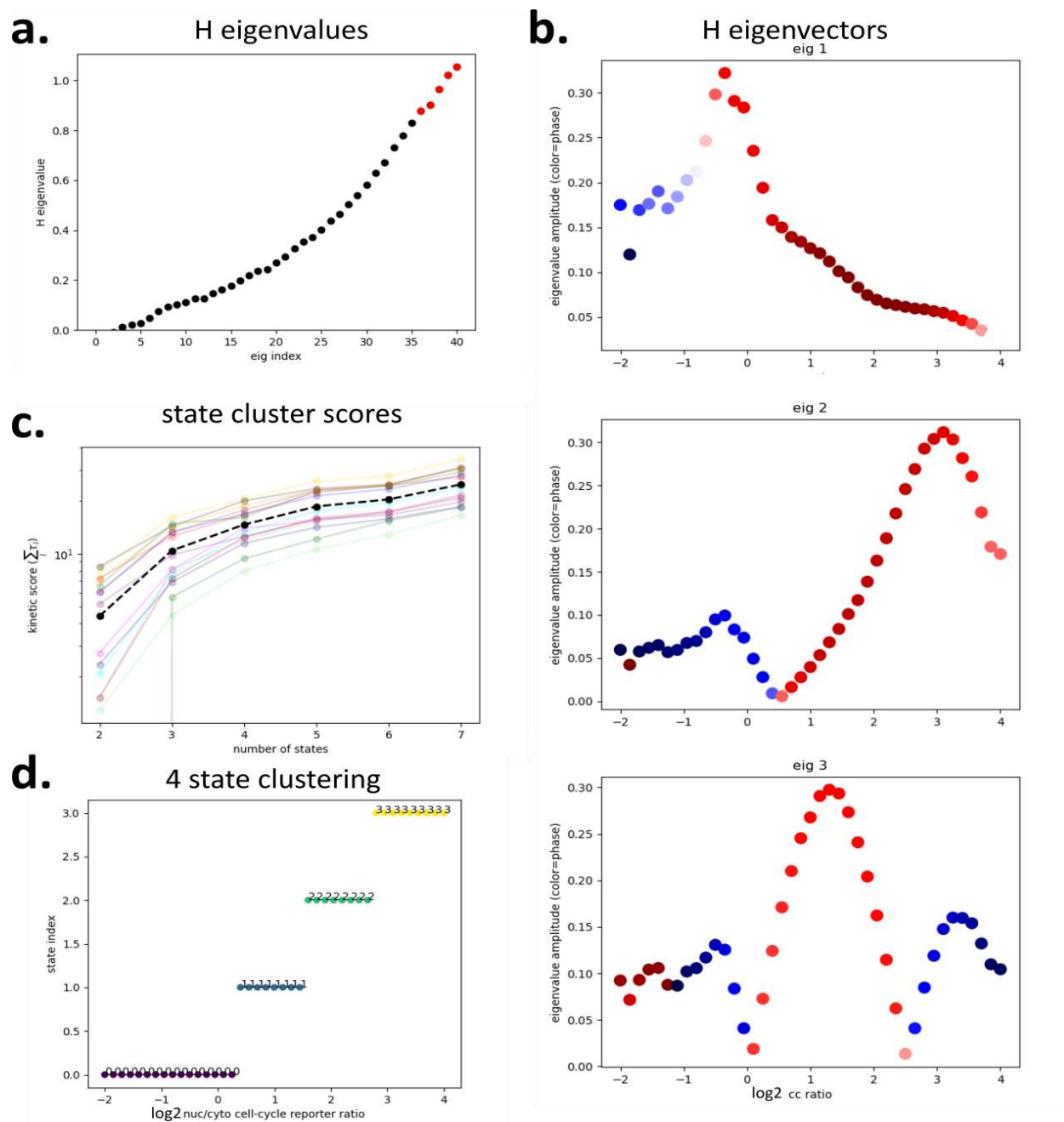
Scatterplots of the first two PCA components for imaging experiments 1 (black) and 2 (red), in overlapping treatments, analyzed in this work. Each dot is a cell before (left) and after (right) applying our batch normalization procedure.



### Supplementary Figure 9: Dynamical clustering of morphodynamical trajectories

- (a) Eigenvalues of the Hermitian extension  $H = \frac{1}{2}[(T + T') + i(T - T')]$  with  $T$  the transition matrix. (b) 2D UMAP of the eigenvectors of  $H$  (each point is a microstate of the transition matrix) colored by absolute value and Euler angle of the complex value. (c) State clustering dynamical information quantified by the sum of the timescales from the eigenvalues of  $T$ . The sum of timescales increases rapidly towards 15 states and begins to saturate. (d) K-means clustering into 14 states.





### Supplementary Figure 10: Dynamical clustering of cell-cycle states

(a) Eigenvalues of the Hermitian extension  $H = \frac{1}{2}[(T + T') + i(T - T')]$  with  $T$  the transition matrix from dividing  $\log_2$  of the cell-cycle reporter levels into 51 microstates. (b) Eigenvectors of  $H$  (each point is a microstate of the transition matrix) colored by Euler angle of the complex value. (c) State clustering dynamical information quantified by the sum of the timescales from the eigenvalues of  $T$ . The sum of timescales increases rapidly towards 4 states and begins to saturate. (d) K-means clustering into 4 states.