

## CELL-CYCLE PHENOTYPING WITH CONDITIONAL RANDOM FIELDS: A CASE STUDY IN *SACCHAROMYCES CEREVISIAE*

Michael B. Mayhew\*      Alexander J. Hartemink<sup>\*,†\*</sup>

\* Program in Computational Biology & Bioinformatics

† Department of Computer Science, Duke University, Durham, NC 27708

### ABSTRACT

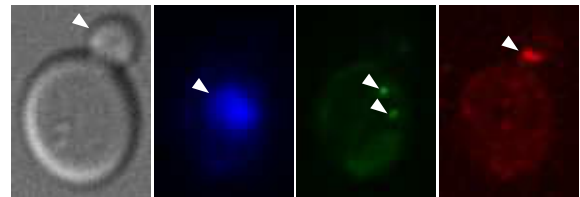
High-resolution, multimodal microscopy grants an intimate view of the inner workings of cells. Complex processes like cell division can be monitored with microscope images, assuming identification of cells and their cell-cycle markers: cellular structures indicative of cell-cycle progress. Here, we explore how spatial relationships between these markers can facilitate their identification. We grew and synchronized *Saccharomyces cerevisiae* cell cultures and then acquired multimodal image data as the cells proceeded through the cell cycle. We trained a conditional random field model to capture pixel-level spatial relationships among three different cell-cycle markers observable in our images. We observed good predictive performance of this pixel-level model on three held-out test images, and performance improved when we used marker-level information from our training data to prune model predictions. Our results support the use of conditional random fields in bioimage labeling and encourage the use of as much multiscale information as available in training data when identifying cell-cycle markers.

**Index Terms**— Cell cycle, budding yeast, conditional random field, multimodal imaging, fluorescence imaging, multiscale information.

### 1. INTRODUCTION

*Saccharomyces cerevisiae*, or budding yeast, has proven an indispensable model organism for the study of numerous complex biological processes. Many such processes, including cell division, can be more finely dissected with microscopy methods [1, 2]. To characterize cell division, traditionally one manually counts a large number of cells (~200–500) from fixed samples collected over time, and records the proportion of cells in each time-point sample that show different markers: structures indicative of cell-cycle progress that can be visualized under the microscope. The classical marker of cell-cycle progression in *Saccharomyces cerevisiae* is the bud, or nascent daughter cell (Figure 1, left).

\*With gratitude for grant funding from the NIH (P50-GM081883-01) and DARPA (HR0011-09-1-0040). Supplementary information is available at <http://www.cs.duke.edu/~amink/publications/>.



**Fig. 1.** A single cell imaged using the four modalities of this study. Cell-cycle markers are indicated by white arrows in each panel. From left to right, modalities (and markers): DIC (bud), blue fluorescence (nucleus), green fluorescence (spindle pole bodies), and red fluorescence (myosin ring).

Fluorescence microscopy grants us access to many other markers of cell-cycle progression such as nuclei, spindle pole bodies (SPBs), and myosin rings (Figure 1, center left, center right, right, respectively), enabling a higher-resolution view of cell division [2].

Many image processing techniques ignore informative spatial relationships between markers that might aid in their identification. For example, since we know that spindle pole bodies are embedded in the nuclear envelope [3], we might expect SPB pixels to occur near nuclear pixels. Conditional random fields (CRFs) are a class of probabilistic graphical models that leverage these kinds of spatial dependencies, as well as rich sets of data-derived features, to represent distributions over label assignments for sequence and image data [4, 5, 6]. CRFs and their extensions have begun to appear in the cell-cycle image analysis community, having been successfully applied to time-lapse microscopy images to detect mitotic events in HeLa cells [7].

In this study, we acquired images of fixed samples of dividing budding yeast cell populations using differential interference contrast (DIC) and fluorescence imaging modalities in order to track three different markers of cell-cycle progression: nuclei, myosin rings, and SPBs. We then trained a CRF model using a combination of both pixel-level and regional intensity features taken from each of the different fluorescence imaging modalities. We demonstrate good performance in predicting markers for three held-out test images with the pixel-level model. We make use of additional marker-level

information in our training data, namely marker size, to improve CRF predictions, and we see further improvement in performance. We close by discussing the importance of multiscale information for cell-cycle marker identification, and its implications for cell counting and cell-cycle phenotyping.

## 2. PIXEL-LEVEL CRF MODEL

### 2.1. Model Details

We assign one of five different categorical labels to each pixel (1 = myosin ring, 2 = SPB, 3 = nucleus, 4 = general intracellular, 5 = noncellular). We implement a simple 4-neighbor CRF (Figure 2) to model local spatial dependence between pixel labels, as well as the dependence of those labels on a collection of image-derived features [8]. We assume the model is spatially invariant: the network topology and parameters are the same everywhere across the image. In a CRF, the probability of any labeling  $\mathbf{Y}$ , given features  $\mathbf{X}$ , is given by the following normalized factor product:

$$P(\mathbf{Y}|\mathbf{X}) = \frac{1}{Z(\mathbf{X})} \prod_{i \in I} \phi(Y_i, X_i) \prod_{i,j \in E} \psi(Y_i, Y_j)$$

where

$$Z(\mathbf{X}) = \sum_{\mathbf{Y}^*} \left[ \prod_{i \in I} \phi(Y_i^*, X_i) \prod_{i,j \in E} \psi(Y_i^*, Y_j^*) \right]$$

is a normalization constant or partition function. The sum in the partition function is over all possible pixel labelings,  $\mathbf{Y}^*$ .

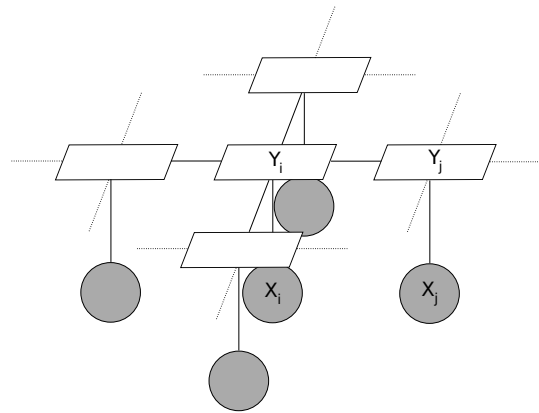
Here,  $I$  is the set of pixels and  $E$  is the set of edges connecting neighboring pixels.  $\mathbf{Y}$  is a particular labeling of a set of pixels,  $\mathbf{X}$  is the set of all features at every pixel, and  $Y_i$  and  $X_i$  are the label and feature set at pixel  $i$ . Our feature set consists of the red, green, and blue fluorescence intensities at each pixel as well as the median green and blue intensities of pixels in three rings (with radii of one, two, and five pixels) centered at each pixel.

The factors in the conditional probability above are also known as node ( $\phi$ ) and edge ( $\psi$ ) potentials, and are defined as follows:

$$\phi(l, f_1, \dots, f_K) = \exp\{\theta_{l,0} + \sum_{k=1}^K \theta_{l,k} f_k\}$$

$$\psi(l, m) = \exp\{\lambda_{l,m,0}\}$$

$\theta$  and  $\lambda$  are parameters of the CRF model. The scope of the node potential includes the label  $l$  and all  $K$  relevant node features. The scope of the edge potential includes the two labels,  $l$  and  $m$ , at each endpoint of an edge. The parameters  $\theta_{l,0}$  and  $\lambda_{l,m,0}$  are bias terms.



**Fig. 2.** Graphical model representation of pixel-level CRF. Our 4-neighbor CRF models the dependence of each label  $Y_i$  on a set of features  $X_i$  as well as on the labels of its four connected neighbors (here,  $Y_j$  and the three other pixels connected to  $Y_i$ ).

### 2.2. Model Training

We built a training data set for the CRF model by hand-labeling 55 cell regions, 55 nuclei, 64 SPBs, and 26 myosin rings in the DIC and fluorescence images of three fields of view. To these pixels, we added pixels from the bounding boxes surrounding each cell region. This procedure resulted in 629,945 total training pixels. We then fit our CRF model by finding the parameters  $\theta$  and  $\lambda$  that maximize the pseudo-likelihood function, an approximation to the full conditional probability above [6, 9]:

$$\arg \max_{\theta, \lambda} \prod_{t=1}^3 \prod_{i \in I_t} \left[ \frac{1}{Z(X_i)} \phi(Y_i, X_i) \prod_{j \in N_i} \psi(Y_i, Y_j) \right]$$

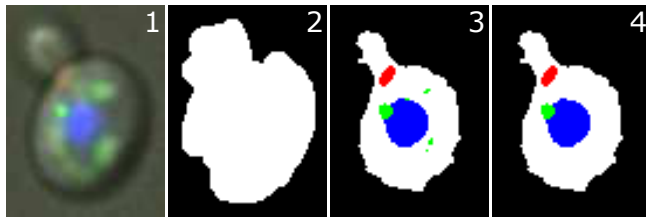
where

$$Z(X_i) = \sum_{Y_i^* \in \{1,2,3,4,5\}} \phi(Y_i^*, X_i) \prod_{j \in N_i^*} \psi(Y_i^*, Y_j)$$

Here, the outer product in the maximization is over the three training images,  $I_t$  is the set of all pixels in training image  $t$ , and  $N_i$  is the set of indices of neighbors of pixel  $i$ . We also trained the model using the sum-product loopy belief propagation algorithm ([10] and references therein), but we did not see much difference in final performance (see Supplement for details). As such, all following results are based on training with the pseudo-likelihood approximation.

### 2.3. Generating Model Predictions

For the sake of computational speed, we generate predictions in a held-out test image only for subregions that potentially include cell and marker pixels. We locate these subregions by using a custom watershed algorithm (Figure 3, center left),



**Fig. 3.** Algorithmic pipeline for pixel labeling. 1) We take as input an image set for a field of view (shown here with registered modalities overlaid). 2) We apply a custom watershed algorithm to identify potential cell regions. 3) We generate a maximum-of-marginals decoding with our CRF model. 4) We remove spurious marker predictions by applying a size filter learned from training data. Myosin ring: red; SPB: green; nucleus: blue; intracellular: white; noncellular: black.

and placing a bounding box around each. We then predict labels for every pixel in the bounding boxes. This resulted in 502,174 pixels to be labeled across all three test images.

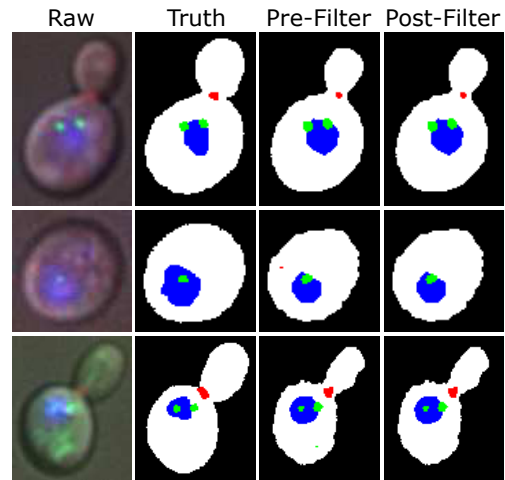
To predict the labels for every pixel, we calculate node and edge potentials based on the features at each pixel and the learned parameters,  $\theta$  and  $\lambda$ . We then use the sum-product loopy belief propagation algorithm to infer the approximate marginal distributions over marker labels at each pixel. We take the label with maximal marginal probability to be the predicted label (Figure 3, center right).

### 3. IMAGE ACQUISITION & PRE-PROCESSING

We grew overnight cultures of two haploid yeast strains in rich media (genotype information in [2]). The populations were synchronized and then released into the cell cycle, with samples collected from the population over time. We fixed the samples and prepared them for imaging. For each time-point sample, we acquired images of at least five fields of view. For each field of view, we acquired a DIC image as well as red, green, and blue fluorescence images. For the green and red fluorescence channels, we imaged z-stacks composed of approximately 20 to 40 slices of  $0.2\mu\text{m}$  thickness; we deconvolved these z-stacks and computed the sum of slices in each deconvolved stack to produce a single green or red fluorescence image. Three fields of view served as training images, and three as test images.

### 4. MODEL PERFORMANCE ON CELL-CYCLE MARKER LABELING

We found that CRF marker predictions were generally biologically sound (Figure 4). Across the three held-out test images, the model often reasonably segmented cell regions from background, predicted one nucleus for each cell, and identified plausible myosin rings and SPBs. However, as seen in



**Fig. 4.** Marker labeling successes. CRF predictions are generally biologically plausible, though in some cases (second and third rows), the size filter improves correspondence between CRF predictions and ground truth. Columns from left to right: overlaid DIC and fluorescent images (raw), hand-labeled cell regions and markers (ground truth), CRF-predicted pixel labels (pre-filter), and model labels after applying size filter (post-filter). Colors are as in Figure 3.

Figures 4 and 5 (third column), the CRF sometimes predicts extra myosin rings and SPBs. To more quantitatively measure predictive performance, we hand-labeled cell-cycle markers in our test images and compared these labelings to predicted labelings from the CRF model. We determined the sensitivity and positive predictive value (PPV) of our CRF model for each marker. As shown in Table 1, the CRF performs well in finding nuclei and SPBs, and more poorly when predicting myosin rings (Table 1). The combination of high sensitivity and lower PPV for myosin rings—and to a lesser extent for SPBs—indicates that while the model overpredicts these markers, the set of model predictions generally includes the set of true markers.

### 5. INCLUSION OF MARKER-LEVEL INFORMATION VIA LEARNED SIZE FILTERS

We noticed that some of the CRF’s extra myosin ring and SPB predictions are small in size compared with the hand-labeled markers in our training data (Figures 4 and 5, third column); some are no more than a few pixels. Our training data not only enables us to learn marker characteristics at the pixel level but also at the (coarser) whole-marker level. So, we decided to incorporate this marker-level information. We computed the minimum pixel area of all myosin rings and SPBs in our training set. We then removed any CRF-predicted myosin ring (or SPB) less than one tenth (or one fifth) the area of the smallest hand-labeled marker of that type (see Supplement for details).

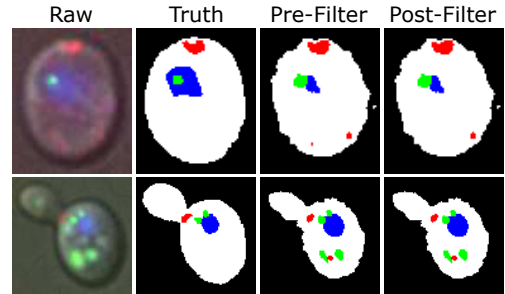
Image <sup>+</sup>	Marker <sup>+</sup>	Sensitivity*	PPV*
SBY1404-t10 (14)	Myosin Ring (2)	1.00 / 1.00	0.08 / 0.14
	SPB (15)	1.00 / 1.00	1.00 / 1.00
	Nucleus (14)	1.00	0.88
SBY1643-t13 (24)	Myosin Ring (15)	0.80 / 0.80	0.57 / 0.63
	SPB (35)	0.91 / 0.86	0.60 / 0.65
	Nucleus (24)	1.00	0.86
SBY1643-t15 (24)	Myosin Ring (14)	1.00 / 0.93	1.00 / 1.00
	SPB (33)	0.97 / 0.97	1.00 / 1.00
	Nucleus (24)	1.00	0.86
Total (62)	Myosin Ring (31)	0.90 / 0.87	0.47 / 0.59
	SPB (83)	0.95 / 0.93	0.79 / 0.83
	Nucleus (62)	1.00	0.86

**Table 1.** <sup>+</sup>Shown in parentheses are true numbers of cell regions and markers. \*Results pre-size filter / post-size filter.

Particularly for myosin ring predictions and slightly for SPB predictions, we found that overall PPV improved (Table 1) with little to no drop in sensitivity. These results are not surprising since coarser-scale information like marker size and shape is not available to the pixel-level CRF. Thus, exploiting available multiscale information appears to be helpful in accurately identifying cell-cycle markers.

## 6. DISCUSSION & FUTURE DIRECTIONS

We have developed a multimodal image acquisition and analysis pipeline to automatically identify markers of cell-cycle progression in images of dividing budding yeast cells. The core of the pipeline is a CRF model that integrates pixel-level features from the different imaging channels and accounts for spatial relationships among markers, at least at the pixel level. Our results demonstrate both the strengths and the limitations of a CRF. While it performs well enough to approximately segment cells from background, and to predict cell-cycle markers reasonably accurately, its discriminative power is limited by the spatial scale of the information at its disposal ([5] and references therein). For example, we may not want to label a patch of pixels as a myosin ring unless we also know marker-level information like the size and shape of the patch, or the distance between the patch and the edge of a cell. Another layer of complexity in marker identification comes from the cell-cycle biology of budding yeast. Specifically, certain numbers and combinations of markers are not plausible in certain budding yeast cells. For example, in the strains we used, a cell cannot contain more than two SPBs, or more than one myosin ring. We are currently investigating model-based approaches to take into account these inter-marker dependencies as well as incorporate more marker-level features, but this depends on precise determination of individual cells (which can be difficult when budded cells are clumped together). In the long term, we will use our size-filtered marker predictions and approximate cell boundaries to automate counting cells with



**Fig. 5.** Marker labeling failures. Some spurious CRF predictions are not removed with a size filter (fourth column). Column layout is the same as in Figure 4. Colors are as in Figures 3 and 4.

different cell-cycle phenotypes. This work is an important step toward that goal, and thus toward more accurate characterization of cell-cycle progression.

## 7. REFERENCES

- [1] S. Di Talia, J.M. Skotheim, J.M. Bean, E.D. Siggia, and F.R. Cross, “The effects of molecular noise and size control on variability in the budding yeast cell cycle,” *Nature*, vol. 448, pp. 947–951, August 2007.
- [2] M.B. Mayhew, J.W. Robinson, B. Jung, S.B. Haase, and A.J. Hartemink, “A generalized model for multi-marker analysis of cell cycle progression in synchrony experiments,” *Bioinformatics*, vol. 27, pp. i295–303, July 2011.
- [3] S.L. Jaspersen and M. Winey, “The budding yeast spindle pole body: Structure, duplication, and function,” *Annu. Rev. Cell & Dev. Biol.*, vol. 20, pp. 1–28, 2004.
- [4] J. Lafferty, A. McCallum, and F. Pereira, “Conditional random fields: Probabilistic models for segmenting and labeling sequence data,” in *International Conference on Machine Learning*, 2001, pp. 282–289.
- [5] X. He, R.S. Zemel, and M.A. Carreira-Perpinan, “Multiscale conditional random fields for image labeling,” in *CVPR*, 2004, pp. 695–702.
- [6] S. Kumar and M. Hebert, “Discriminative random fields: A discriminative framework for contextual interaction in classification,” in *ICCV*, 2003, pp. 1150–1157.
- [7] L. Liang, X. Zhou, F. Li, S.T.C. Wong, J. Huckins, and R.W. King, “Mitosis cell identification with conditional random fields,” in *IEEE/NIH Life Science Systems and Applications Workshop*, 2007, pp. 9–12.
- [8] M. Schmidt, “UGM: Matlab code for undirected graphical models,” <http://www.di.ens.fr/~mschmidt/Software/UGM.html>.
- [9] J. Besag, “Statistical analysis of non-lattice data,” *Journal of the Royal Statistical Society, Series D (The Statistician)*, vol. 23, pp. 179–195, September 1975.
- [10] K.P. Murphy, Y. Weiss, and M.I. Jordan, “Loopy belief propagation for approximate inference: An empirical study,” in *Proceedings of Uncertainty in AI*, 1999, pp. 467–475.