

Automated detection and quantification of single RNAs at cellular resolution in zebrafish embryos

L. Carine Stapel, Benoit Lombardot, Coleman Broaddus, Dagmar Kainmueller, Florian Jug, Eugene W. Myers and Nadine L. Vastenhouw*

ABSTRACT

Analysis of differential gene expression is crucial for the study of cell fate and behavior during embryonic development. However, automated methods for the sensitive detection and quantification of RNAs at cellular resolution in embryos are lacking. With the advent of single-molecule fluorescence *in situ* hybridization (smFISH), gene expression can be analyzed at single-molecule resolution. However, the limited availability of protocols for smFISH in embryos and the lack of efficient image analysis pipelines have hampered quantification at the (sub)cellular level in complex samples such as tissues and embryos. Here, we present a protocol for smFISH on zebrafish embryo sections in combination with an image analysis pipeline for automated transcript detection and cell segmentation. We use this strategy to quantify gene expression differences between different cell types and identify differences in subcellular transcript localization between genes. The combination of our smFISH protocol and custom-made, freely available, analysis pipeline will enable researchers to fully exploit the benefits of quantitative transcript analysis at cellular and subcellular resolution in tissues and embryos.

KEY WORDS: Cell segmentation, Gene expression, Image analysis pipeline, smFISH, Zebrafish

INTRODUCTION

Analysis of gene expression patterns is an essential tool in many areas of biological research. In developmental biology, for instance, it provides valuable information on the role of differential gene expression in determining cell fates (Junker et al., 2014a; Satija et al., 2015; Thisse and Thisse, 2008; Tomancak et al., 2007). Spatial patterns of gene expression have historically been studied by RNA *in situ* hybridization, but this technique is generally not quantitative (Gross-Thebing et al., 2014; Thisse and Thisse, 2008; Tomancak et al., 2007). Relative levels of gene expression are often studied by RNA-sequencing approaches. When performed at the cellular level, however, this technique only detects the ~10% most abundant transcripts and is thus rather insensitive (Grün et al., 2014; Junker et al., 2014a; Satija et al., 2015). Furthermore, neither technique provides subcellular resolution. The development of single-molecule fluorescence *in situ* hybridization (smFISH) has enabled the detection of individual transcripts both in single cells and tissues (Bahar Halpern et al., 2015; Battich et al., 2013; Boettiger and Levine, 2013; Itzkovitz et al., 2012, 2011; Little et al., 2013; Lyubimova et al., 2013; Mueller et al., 2013; Nair et al., 2013; Oka and Sato, 2015; Peterson et al., 2012; Raj et al., 2008). This

technical advance has, for example, improved our understanding of the design principles of the developing mouse intestine (Itzkovitz et al., 2012) and the establishment of precise developmental gene expression patterns in *Drosophila* blastoderm embryos (Boettiger and Levine, 2013; Little et al., 2013). However, broad application of smFISH in complex samples has been hampered by the limited availability of protocols for embryos and by the lack of an automated image analysis pipeline that combines transcript detection with cell segmentation (Bahar Halpern et al., 2015; Itzkovitz et al., 2011; Lyubimova et al., 2013; Oka and Sato, 2015). Thus, the potential of smFISH in fields such as developmental biology remains to be fully exploited.

Here, we present a protocol for smFISH on embryo sections in combination with an analysis pipeline for automated transcript detection and cell segmentation. We apply our approach to the quantification of RNA expression in single cells of developing zebrafish embryos. To illustrate the power of our method, we identified cell type-specific differences in gene expression and assigned transcripts to different subcellular compartments. The combination of our smFISH protocol and image analysis pipeline opens the door for automated, high-resolution transcript analysis in a variety of complex systems. This tool will be valuable in many areas of biological research, including development, stem cell biology and regeneration.

RESULTS AND DISCUSSION

Sensitive and specific detection and quantification of transcripts

To detect mRNA at single-molecule resolution, we developed a protocol for smFISH on 8 µm cryosections of zebrafish embryos. We imaged and analyzed stacks of 17 z-slices with 0.3 µm spacing, corresponding to a total thickness of ~5 µm (Fig. 1A and Materials and Methods). To visualize single RNA molecules, we used 48 oligonucleotide probes 20 bases long, each coupled to one fluorophore (Stellaris, Biosearch Technologies) (Raj et al., 2008). Once hybridized to an RNA molecule, the probes generate diffraction-limited fluorescent spots that can readily be distinguished from background signal (Fig. 1).

To test our protocol, we performed smFISH for *ntla* (also known as *ta* - ZFIN) and *eif4g2a* on sections of zebrafish embryos at 50% epiboly [5.3 hours post fertilization (hpf)] (Fig. 1B–E, Fig. S1). *ntla* is involved in mesoderm specification and has been shown to be expressed in the presumptive mesoderm at the margin of the embryo (Harvey et al., 2010; Schier and Talbot, 2005) (Fig. S2A,B). By contrast, *eif4g2a* is a ubiquitously expressed housekeeping gene (Fig. S2C). To detect transcripts for both genes simultaneously, we labeled the two probe sets with different fluorophores (*ntla*-Q670, *eif4g2a*-CF610). We included DAPI staining to detect nuclei (Fig. 1D,E). Embryos were imaged in a tile scan on a wide-field microscope and the resulting images were stitched with the *Grid/*

Max Planck Institute of Molecular Cell Biology and Genetics, Pfotenhauerstr. 108, Dresden 01307, Germany.

*Author for correspondence (vastenhouw@mpi-cbg.de)

Received 24 July 2015; Accepted 14 December 2015

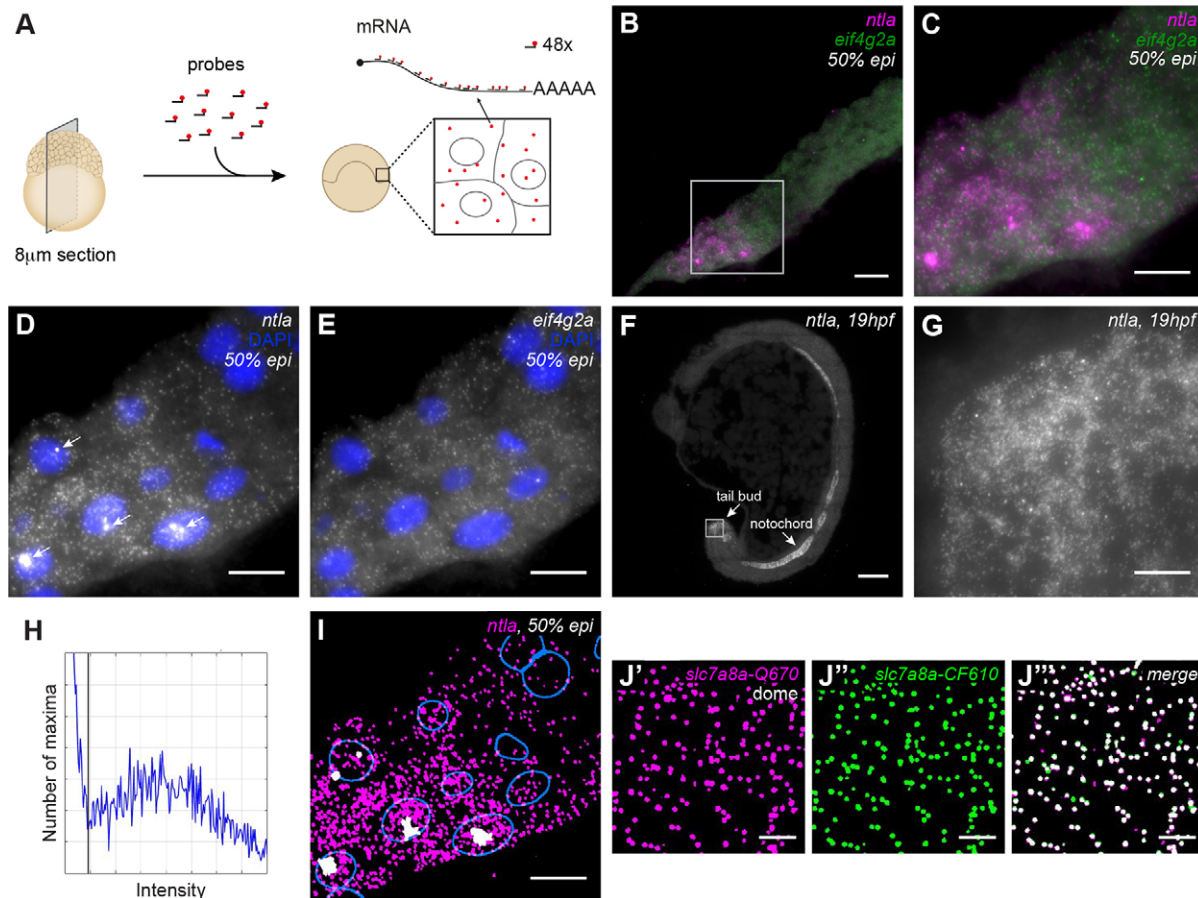


Fig. 1. Sensitive and specific detection of transcripts and transcription foci in zebrafish sections with smFISH. (A) Overview of smFISH method on sections of zebrafish embryos. (B) smFISH for *ntlA* and *eif4g2a* on a cryosection of a 50% epiboly stage embryo. Scale bar: 20 μ m. Complete animal cap in Fig. S1A. (C-E) Box in B at higher magnification. Scale bars: 10 μ m. (C) Dual-color smFISH for *ntlA* (magenta) and *eif4g2a* (green). (D) smFISH for *ntlA*. Arrows indicate transcription foci. (E) smFISH for *eif4g2a*. (F) smFISH for *ntlA* on a 19 hpf embryo. Scale bar: 100 μ m. (G) Detail of F, showing smFISH for *ntlA* in the tail bud. Scale bar: 10 μ m. (H) Spot intensity plot for *ntlA* smFISH on the complete animal cap shown in Fig. S1A. Black line indicates transcript detection threshold. (I) Detection of *ntlA* transcripts (magenta) and foci (white) with the Transcript analysis plugin. Nuclear outlines are indicated in blue. Scale bar: 20 μ m. (J) Dual-color detection of *slc7a8a* with two non-overlapping, differently labeled probe sets. Scale bars: 5 μ m. (J') smFISH *slc7a8a*-Quasar670. (J'') smFISH *slc7a8a*-CalFluor610. (J''') Dual-color view of transcripts detected with the two non-overlapping probe sets in J' and J''. Images are maximum projections of 17 z-slices spaced by 0.3 μ m.

Collection stitching plugin in Fiji (Preibisch et al., 2009; Schindelin et al., 2012). In agreement with its known expression pattern, *ntlA* expression was only detected at the margin of the embryo (Fig. 1B-D, Fig. S1). By contrast, *eif4g2a* was detected ubiquitously (Fig. 1B,C,E, Fig. S1A). Interestingly, and consistent with the localization of the upstream activators of *ntlA* in the yolk syncytial layer [BMP and Nodal (Harvey and Smith, 2009; Harvey et al., 2010; Schier and Talbot, 2005)], smFISH revealed that there is a vegetal-animal gradient of *ntlA* expression (Fig. 1B-D, Fig. S1). *ntlA* was also detected at single-molecule resolution in notochord and tail bud at 19 hpf (Fig. 1F,G), in line with whole-mount *in situ* hybridization data (Schier and Talbot, 2005), illustrating the versatility of our protocol. Taken together, these results indicate that we can obtain specific, high-resolution information on gene expression for multiple genes simultaneously in zebrafish embryos at various stages of development.

Next, we developed a Fiji plugin (*Transcript analysis*) to quantify transcript numbers in an automated fashion. To detect transcripts, we filtered images, detected local maxima of intensity and used a threshold to separate true transcripts from background noise, similar to previous approaches (Lyubimova et al., 2013; Mueller et al.,

2013; Raj et al., 2008). To determine the appropriate threshold for detection of *ntlA* transcripts, we plotted the intensity distribution of all detected maxima (Fig. 1H). For each probe set, we manually set the threshold for transcript detection between the low intensity peak, reflecting background signal, and the high intensity peak, reflecting transcripts. The unimodal shape of the transcript peak confirms that the spots we identify were indeed single RNA molecules (Raj et al., 2008; Vargas et al., 2005). Comparison of the transcript detection output with the smFISH image suggested that the sensitivity of transcript detection with the image analysis pipeline is high (Fig. 1D,I).

To quantify the sensitivity and specificity of our method, we first analyzed *slc7a8a* transcripts with two probe sets that were labeled with different fluorophores (Fig. 1J). Of the spots detected with probe set 1 (*slc7a8a*-Quasar670), 87% was also detected with probe set 2 (*slc7a8a*-CalFluor610). Conversely, 81% of spots detected with probe set 2 was also detected with probe set 1. This might even be an underestimation of the efficiency, because the use of two probe sets for one gene precludes the use of the 48 best probes. In comparison, previous studies reported detection efficiencies of 70-85% for smFISH (Oka and Sato, 2015; Raj et al., 2008). Next, to test

the specificity of the method, we performed dual-color labeling of two different genes (*eif4g2a* and *ntla*). This resulted in an overlap of only 2% in cells where both genes are expressed (Fig. S3). Finally, transcript numbers obtained by smFISH correlated well ($r=0.94$) with RNA-sequencing data (Pauli et al., 2012), confirming the quantitative power of our smFISH approach (Fig. S4). Taken together, these results show that our method detects transcripts efficiently and specifically.

In addition to individual transcripts, high-intensity foci corresponding to sites of active transcription (Bahar Halpern et al., 2015; Levesque and Raj, 2013) were sometimes observed in the nucleus (Fig. 1D, arrows). As expected, a maximum of two foci per nucleus was observed, one for each allele. We extended our analysis pipeline to include the automated detection of transcription foci based on their size and intensity (Materials and Methods and Fig. 1I). We compared detected foci with foci in smFISH images and found a detection sensitivity close to 90%, with a precision of more than 97% (Fig. S5). Only weak foci were not detected automatically. When 100% detection efficiency of foci is essential, an intronic probe can be used to mark transcription sites specifically. To quantify the number of transcripts in each focus, we divided the sum intensity of the transcription foci by the median sum intensity of the transcripts (Mueller et al., 2013). In conclusion, our smFISH protocol and analysis pipeline (Fig. S6) enable the detection of single RNA molecules and transcription foci in zebrafish embryo sections with high sensitivity and specificity.

An automated membrane segmentation pipeline to assign transcripts to cells

In order to assign transcripts to cells and specific cellular compartments, cells and nuclei have to be segmented. So far, the use of smFISH for the quantitative analysis of gene expression in complex samples has been hampered by the lack of an efficient cell segmentation pipeline. Current analysis pipelines rely on manual segmentation of cells (Bahar Halpern et al., 2015; Itzkovitz et al., 2011; Lyubimova et al., 2013; Oka and Sato, 2015), which is not feasible for large amounts of data or samples as large as the zebrafish embryo. To overcome this problem, we developed an automated pipeline to segment cells in tissue sections (Fig. 2, Fig. S6).

To identify the cell membrane, we incorporated a phalloidin-staining step in our smFISH protocol (Fig. 2A). We used the middle slice of our z -scan acquisition for cell segmentation. This is a good approximation of the cell outline in thin sections. We trained a cascaded Random Forest (Breiman, 2001; Tu and Bai, 2010) to predict for each pixel the probability that it belongs to the membrane, and additionally the probability that it belongs to a membrane intersection point (vertex) based on the phalloidin staining (Fig. 2B). Given these probabilities, we can trace paths that are likely to run along the membrane between points that are most likely vertices. This results in a mask of cell membranes (Fig. 2C). Depending on the quality of the membrane staining, the membrane-tracing software can produce both over- and under-segmentation errors. These errors can easily be corrected manually by drawing missing lines and breaking excessive ones with our Fiji tool 'Cell annotation'. In our samples, and with the settings we chose, automated segmentations exhibit on average 91% precision (100% would indicate no over-segmentation) and 70% recall (indicating the fraction of correct segmentations) (Fig. S7). Manual corrections take 5 min per image, compared with 20 min for a completely manual segmentation. Finally, the individual cells are identified (Fig. 2D). Our pipeline significantly reduces cell segmentation time

compared with existing approaches that rely on manual segmentation (Bahar Halpern et al., 2015; Itzkovitz et al., 2011; Lyubimova et al., 2013; Oka and Sato, 2015). In the future, it might be possible to implement assisted manual correction, which would further reduce segmentation times. In addition, we segmented nuclei to be able to distinguish between cytoplasmic and nuclear transcripts (Fig. 2E,F). We used a watershed-based approach (Ollion et al., 2013) to segment nuclei in 2D on a maximum z -projection. Together, our smFISH method, cell segmentation and nuclear detection allow us to automatically assign transcripts and transcription foci to specific cells and nuclei (Fig. 2F).

Using the automated pipeline, we can calculate transcript densities per cell as number of transcripts per μm^3 . We used transcript density as a measure of gene expression because it has been shown to be a more reliable readout than transcript number (Padovan-Merhar et al., 2015), and because we do not image complete cells in our cryosections. A flowchart of the complete analysis pipeline including transcript detection can be found in Fig. S6.

Quantification of cell type-specific differences in gene expression

To validate our method, we quantified gene expression at dome stage (4.3 hpf) when the first two cell types, the extra-embryonic cells of the enveloping layer (EVL) and the embryonic cells (deep layer, DEL) (Kimmel et al., 1990), have been specified (Fig. 3A-C, Fig. S8). We analyzed the maternally loaded gene *eif4g2a* and two genes involved in early zebrafish development, *sox19a* and *mex3b*. No differences in gene expression were detected for these genes by regular *in situ* hybridization (Fig. S9). To quantify gene expression in EVL and DEL, we expanded our annotation tool to categorize cells. With this tool, any segmented cell can be assigned to a selected class by simply clicking on it (Fig. 3D). Here, we identified cells based on location, but markers can also be used. Antibody staining can easily be incorporated in the smFISH protocol (data not shown; Lyubimova et al., 2013; Raj et al., 2008). Quantification of transcript densities in EVL and DEL revealed that expression of *sox19a* was 4.6-fold higher in the EVL than in the DEL, whereas expression of *mex3b* was 5.1-fold lower in the EVL (Fig. 3E). By contrast, *eif4g2a* was expressed at similar levels in both cell types (Fig. 3E). Thus, our approach allows sensitive detection and quantification of differences in gene expression between cells in an embryo, making it a useful tool in a variety of applications, such as the analysis of transcript levels in relation to cell fate determination.

Quantification of subcellular transcript distribution

The localization of mRNAs plays an important role in organizing cellular function (Besse and Ephrussi, 2008; Jambor et al., 2015; Lécuyer et al., 2007). To determine whether our approach is able to identify differences in mRNA localization, we assigned transcripts of three genes to nuclei and cytoplasm and identified the level of transcriptional activity (in transcription foci) at sphere stage (4 hpf) (Fig. 4A-C, Fig. S9). The maternally loaded housekeeping gene *eif4g2a* was expressed at an average density of 8.1×10^{-2} transcripts per μm^3 . Very few transcripts were found in foci or dispersed throughout the nucleus and most *eif4g2a* transcripts were localized to the cytoplasm (Fig. 4A,D, Fig. S10A). Thus, at sphere stage, most *eif4g2a* transcripts are available for translation. The zygotically expressed genes *tbx16* (spadetail) and *akap12b* were expressed at average densities of 3.0×10^{-2} and 4.2×10^{-2} transcripts per μm^3 , respectively. In contrast to *eif4g2a*, a large proportion of these transcripts was located in transcription foci or scattered throughout

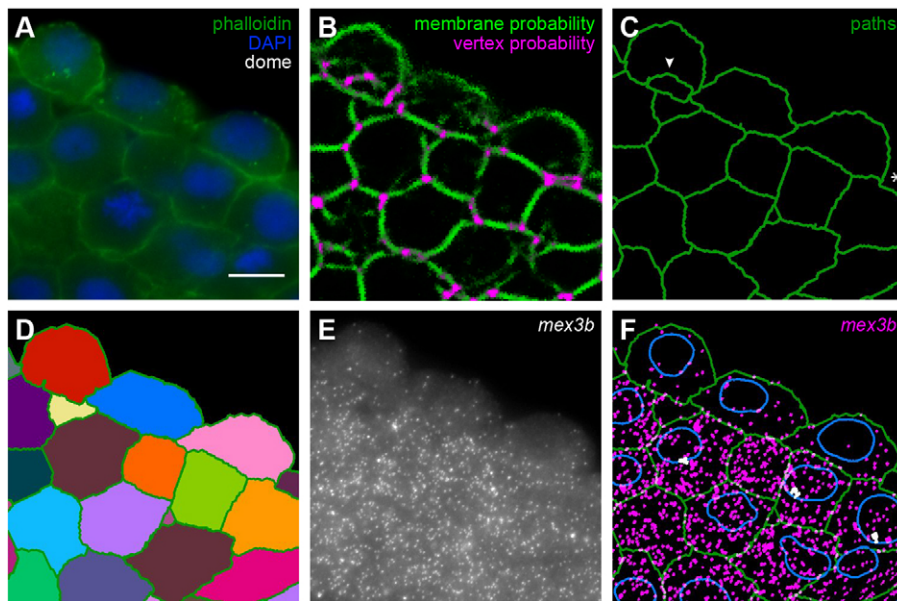


Fig. 2. Automated membrane detection to assign transcripts to cells and nuclei.

(A) Phalloidin staining (green) and DAPI staining (blue) on an smFISH sample to identify cell membrane and nuclei. (B) Output of the Cascaded Random Forest classification for membrane (green) and vertex points (magenta) performed on A. (C) Membrane traces (green) and the vertex points from B as input. (D) Cell mask after manual correction of membrane traces. (E) smFISH for *mex3b*. (F) Combined detection of transcripts (magenta), transcription foci (white), outlines of cells (green) and nuclei (blue). Scale bar: 10 μm .

the nucleus (Fig. 4B-D, Fig. S10B,C). Fewer than half of *tbx16* and *akap12b* transcripts were located in the cytoplasm (Fig. 4B-D, Fig. S10B,C). Since nuclei were segmented in 2D, small nuclear sizes might reflect incomplete presence of the nucleus in the *z*-stack, resulting in the mis-assignment of transcripts. To avoid this potential problem, we analyzed only those cells with the top 25%

largest nuclei, which are most likely to fill the entire *z*-stack (Fig. 4D, Fig. S11A). However, analyzing all cells resulted in very similar distributions (supplementary Materials and Methods and Fig. S11B). Taken together, these data show that our approach can quantify the distribution of transcripts between nuclei and cytoplasm.

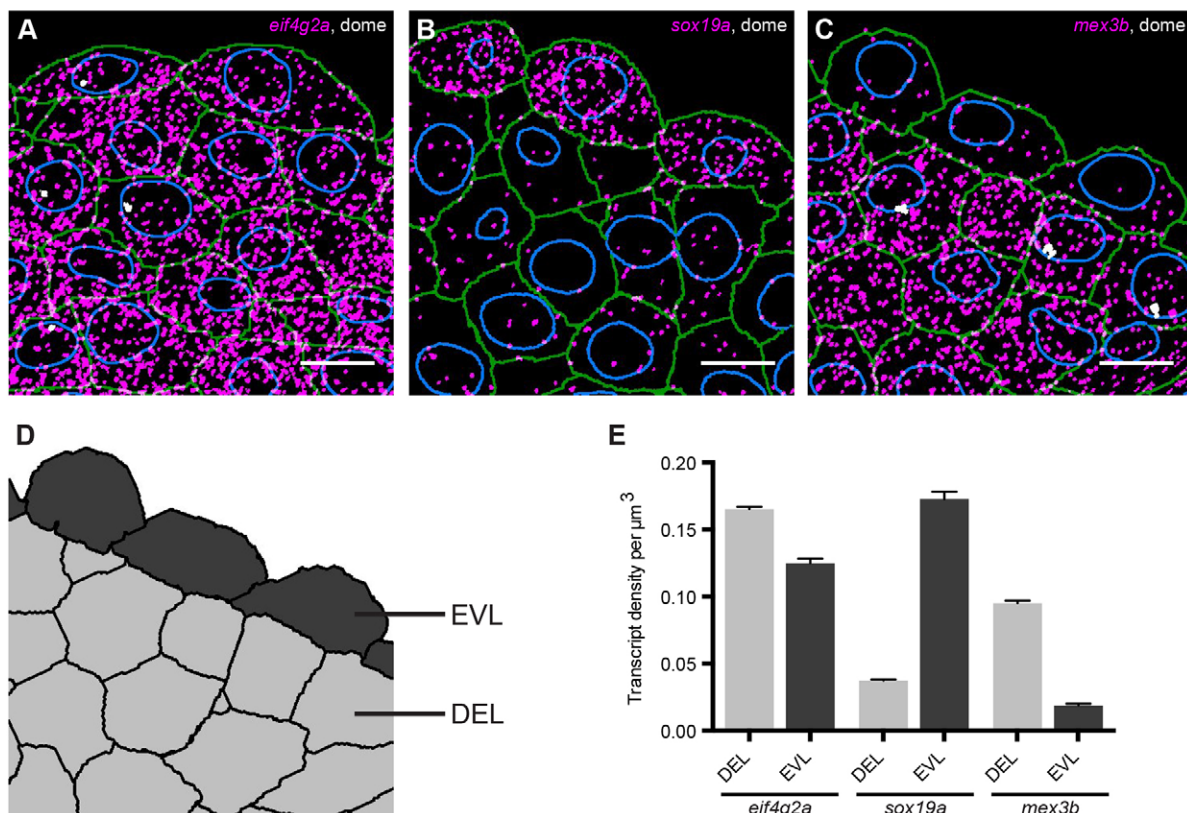


Fig. 3. smFISH provides quantitative spatial information on gene expression. Detected transcripts (magenta) and transcription foci (white) for *sox19a* (A), *mex3b* (B) and *eif4g2a* (C) at dome stage. Green, cell outlines; blue, nuclear outlines. Scale bars: 10 μm . Images are maximum projections of 17 *z*-slices spaced by 0.3 μm . (D) Tissue mask corresponding to C, to distinguish between EVL and DEL cells. (E) Quantification of transcript levels in DEL and EVL. Values are means from sections of three different embryos. Error bars represent s.e.m.

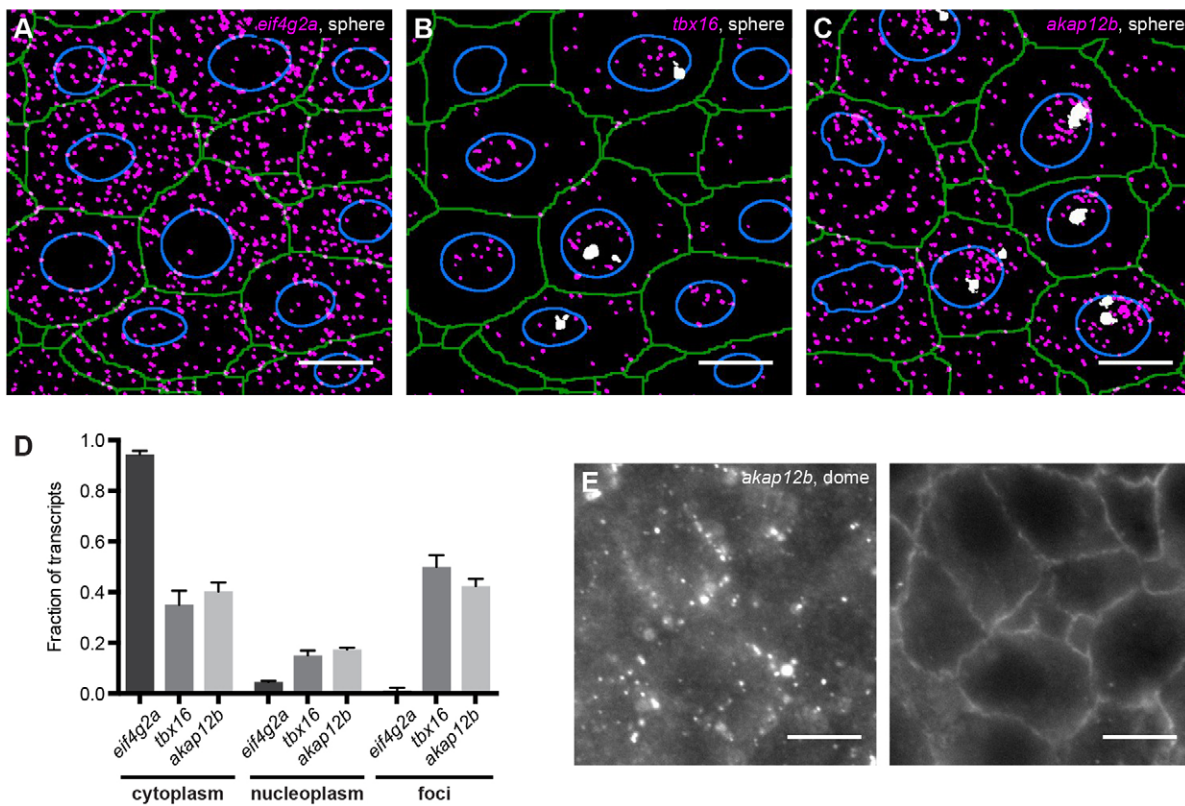


Fig. 4. smFISH provides quantitative subcellular information on gene expression. Detected transcripts (magenta) and transcription foci (white) for *eif4g2a* (A), *tbx16* (B) and *akap12b* (C) at sphere stage. Green, cell outlines; blue, nuclear outlines. Scale bars: 10 μ m. Images are maximum projections of 17 z-slices spaced by 0.3 μ m. (D) Single transcript quantification. Values are averages from sections of three different embryos. Error bars represent s.e.m. (E) smFISH for *akap12b* (left) with corresponding membrane staining for Phalloidin (right) at dome stage. Images are single z-slices. Scale bars: 10 μ m.

mRNAs can be localized to more sites than nuclei and cytoplasm (Jambor et al., 2015; Lécuyer et al., 2007). Interestingly, and in contrast to the localization of *akap12b* at sphere stage, at the onset of gastrulation, most *akap12b* transcripts were localized in clusters at the plasma membrane (Fig. 4E). *akap12b* encodes a scaffold protein that regulates the transition from convergence to extension movements during gastrulation (Weiser et al., 2007). The zebrafish Akap12b protein has been shown to localize to plasma membranes when expressed in cultured human cells, but not much was known about the potential localization of *akap12b* mRNA (Weiser et al., 2007). Localization of *akap12b* mRNA to the membrane might facilitate its translation right at the site of action of the protein (Besse and Ephrussi, 2008; Lécuyer et al., 2007). Taken together, these results show that our approach can quantify asymmetries in the localization of transcripts, which is important for determining their function.

In conclusion, we have developed a method in zebrafish that enables the automated detection and quantification of transcripts at cellular and subcellular resolution in large samples. So far, studies in large and complex samples have used manual segmentation to assign transcripts to specific cells (Bahar Halpern et al., 2015; Itzkovitz et al., 2011; Lyubimova et al., 2013; Oka and Sato, 2015). This has limited the number of cells that could be analyzed, and as a consequence, the potential of smFISH has not been fully exploited. For example, to draw reliable conclusions about variability in gene expression between cells, data on large numbers of cells is required (Battich et al., 2013). Furthermore, gene expression has often been indicated as a function of an animal/organ axis (Hoyle and Ish-Horowicz, 2013; Junker et al., 2014b; Kim et al., 2013; Nair et al.,

2013; Peterson et al., 2012). Although this kind of representation is informative, cellular resolution would provide more precise information. Recent examples of where this would be of value include sonic hedgehog signaling dynamics in the developing neural tube (Peterson et al., 2012) and the relationship between the expression level of a micro RNA and its target (Kim et al., 2013). In summary, our method facilitates the automated detection and quantification of transcripts and their assignment to cells and subcellular structures. Our custom-made software is freely available in KNIME and Fiji and allows researchers working with complex tissues in diverse systems to start exploiting the benefits of high-resolution transcript analysis.

MATERIALS AND METHODS

Zebrafish

Zebrafish were maintained and raised under standard conditions. Wild-type (TLAB) embryos were left to develop to the desired stage at 28°C. Staging was done according to Kimmel et al. (1995).

smFISH

smFISH sample preparation

Embryos were fixed in 4% formaldehyde in PBT (PBS with 0.1% Tween) at 4°C overnight. The next day, embryos were dechorionated manually in PBT and incubated in several changes of fresh 30% sucrose in PBS over the course of several hours before being incubated in 30% sucrose in PBS/OCT (50/50, v/v) at 4°C for 5 days. Then, embryos were embedded in OCT and blocks were quickly frozen in precooled isopentane at –80°C. Cryosection blocks were wrapped in foil and stored at –80°C. 8 μ m cryosections were attached to selected #1.5 22×22 mm coverslips, that were cleaned by sonicating once in 1:20 mucasol and twice in 100% ethanol, and were then

coated with 1:10 poly-L-lysine (Sigma, P8920). Coverslips with sections were stored in sealed 6-well plates at -80°C .

smFISH

smFISH was performed as described previously (Lyubimova et al., 2013) with some changes to obtain high-quality sections of fragile embryos and to reduce background signal. In brief, sections were postfixed in 4% paraformaldehyde in PBS for 15 min and rinsed twice with PBS. Sections were equilibrated in 70% ice-cold ethanol for 5 min and incubated in fresh 70% ice-cold ethanol at 4°C for 4-8 h for permeabilization. Samples were rehydrated in $2\times$ SSC and subjected to a mild proteinase K digestion step at 1:2000 (10 mg/ml stock) for 10 min to increase accessibility of RNAs. After two 5 min washes in $2\times$ SSC, samples were equilibrated in 10% smFISH wash buffer for several minutes (10% formamide, $2\times$ SSC) before probe hybridization. Probes (Biosearch Technologies) were hybridized at a concentration of 75-250 nM in 10% hybridization buffer [10% dextran sulfate (w/v) (Sigma, D8906), 10% formamide (v/v), 1 mg/ml *E. coli* tRNA (Roche), 0.02% BSA, 2 mM vanadyl-ribonucleoside complex (NEB, S1402S), $2\times$ SSC]. For this, smFISH wash buffer was carefully drained from the coverslips as much as possible before coverslips were placed section down on a 100 μl drop of hybridization buffer with probe on a Parafilm-coated cell culture dish. Hybridization was performed at 30°C for ~ 16 h. Then, coverslips were carefully released from the Parafilm with 10% wash buffer. Samples were rinsed with 2 ml of 10% smFISH wash buffer and washed for 2×30 min with 1 ml 10% wash buffer at 30°C . 1:2500 DAPI (1 mg/ml stock) and 1:100 Phalloidin (Life Technologies, A12379) were added to the second wash to stain the nucleus and membrane, respectively. After the second wash, samples were placed in GLOX buffer (10 mM Tris-HCl, pH 7.5, 0.4% glucose, $2\times$ SSC) at 4°C until mounting. Samples were mounted in freshly prepared GLOX mounting medium [GLOX buffer with 1:50 each of 3.7 mg/ml glucose oxidase (Sigma, G2133), Catalase suspension (Sigma, C3515) and Trolox] and sealed with nail polish.

smFISH probes

A total of 48 probes per mRNA, each 20 bases long, were designed using the Stellaris Probe Designer (<https://www.biosearchtech.com/stellarisdesigner/>). CAL Fluor Red 610 and Quasar 670 labeled probes were ordered from Biosearch Technologies. For *slc7a8a*, we designed 96 probes and ordered them with alternating fluorophores for dual-color detection. For probe sequences, see Table S1.

smFISH imaging

Samples were imaged in a tile scan of 19 *z*-sections on a Delta Vision epifluorescence microscope equipped with a 60×1.42 NA oil objective, a Photometrics Cool Snap CCD camera and the following emission filter sets: 435/48, DAPI; 525/36, Alexa Fluor 488; 632/60, CAL Fluor Red 610; 676/34, Quasar 670. Pixel size in the image plane is 0.1072×0.1072 μm . We acquired *z*-stacks with 0.3 μm spacing. After acquisition, image tiles were stitched with the 'Grid/Collection stitching' plugin in Fiji (Preibisch et al., 2009; Schindelin et al., 2012).

Image analysis

The first 17 optical *z*-slices (corresponding to ~ 5 μm thickness) of our 8 μm sections were used for analysis. We empirically determined that this gives the best smFISH results. For other tissues and probe sets, the depth at which good imaging results can be obtained with an epifluorescence microscope might differ depending on the overall background levels (auto fluorescence) and non-specific probe binding. Therefore, when setting up the technique in another tissue, the thickness of sections and the imaging depth should be empirically determined.

Transcript detection

First, background signal was removed from images using top-hat filtering. Next, images were smoothed with a Gaussian kernel to remove noise. Transcripts were detected as local maxima in this image and distinguished from the background noise with an intensity threshold, T_{tx} . In the histogram of

local maxima intensity, T_{tx} was chosen between the one or two sharp peak(s) corresponding to the background and the lower peak of the transcripts at higher intensity. Transcripts were segmented using watershed segmentation initiating from the detected maxima. Transcription foci were detected among the regions defined in the transcript segmentation with the use of thresholds for maximum intensity and volume. For further details, see supplementary Materials and Methods.

Cell segmentation

Cell segmentation was based on Phalloidin staining. The middle slice of the *z*-scan acquisition was used for cell segmentation as an approximation of the cell outline in our thin sections. With a pixel-level classifier, the probability of being on a membrane, as well as a probability of being at the intersection of multiple membranes (i.e. a vertex) was assigned to each pixel. To this end, we trained a two-level cascaded Random Forest classifier from manually segmented training data. Based on the output of this classifier, we traced membranes as highly likely paths between vertices. The set of shortest paths whose length falls below a specific threshold constitutes our automated membrane segmentation. For more details, see supplementary Materials and Methods.

Imaging software

The most recent version of the image analysis software described in this paper, as well as the documentation, is available via <http://tinyurl.com/KNIME-MS-ECS> and <http://fiji.sc/MS-ECS-2D>.

Whole-mount *in situ* hybridization

Whole-mount *in situ* hybridization was performed as described previously (Thisse and Thisse, 2008). After staining, embryos were cleared in methanol and gradually transferred to 87% glycerol for imaging. Samples were imaged in 87% glycerol on a Leica M165C dissecting scope equipped with a Leica MC170 HD camera. Probes were made by PCR amplification of regions of target gene cDNA and cloning these into the pSC-A vector (StrataClone PCR cloning kit). The following primer pairs were used: *eif4g2a* FW: ACGCTTCTCTTTGGCCTCATCG, RV: CAGGCTGTGT-TTGGTAATCCCTG; *sox19a* FW: GAATGACCCAGCTGAACGGTGG, RV: GCCATGGCGGATGGATACTGC; *mex3b* FW: CCCTGCGAGCA-AAGACCAATAC, RV: CGTCCCCATGCAGGTCAAAC. For *ntla*, a previously published probe was used (Bennett et al., 2007).

Acknowledgements

We thank Caren Norden, Máté Pálffy, Iain Patten and Pavel Tomancak for critically reading the manuscript, Jan Philipp Junker and Alexander van Oudenaarden for advice on smFISH, and the following MPI-CBG Services and Facilities for their support: Biomedical Services (Fish facility), Light Microscopy, and Scientific Computing.

Competing interests

The authors declare no competing or financial interests.

Author contributions

L.C.S. and N.L.V. conceived the study. L.C.S. performed and analyzed all experiments. B.L. developed and documented the transcript detection, nuclear segmentation and tissue annotation pipelines. C.B., D.K. and E.W.M. developed and documented the membrane segmentation pipeline. F.J. implemented the membrane and vertex classification tools in KNIME. L.C.S. and N.L.V. wrote the manuscript with input from other authors.

Funding

This work was supported by Max Planck Institute of Molecular Cell Biology and Genetics core funding; the German Federal Ministry of Research and Education [031A099]; a Human Frontier Science Program Career Development Award [CDA-00060/2012-C to N.L.V.]; the Klaus Tschira Stiftung (to E.W.M.); a Boehringer Ingelheim Fonds PhD fellowship (to L.C.S.); and a Center for Systems Biology Dresden postdoctoral fellowship (to D.K.).

Supplementary information

Supplementary information available online at <http://dev.biologists.org/lookup/suppl/doi:10.1242/dev.128918/-DC1>

References

Bahar Halpern, K., Tanami, S., Landen, S., Chapal, M., Szliak, L., Hutzler, A., Nizhberg, A. and Itzkovitz, S. (2015). Bursty gene expression in the intact mammalian liver. *Mol. Cell* **58**, 147-156.

- Battich, N., Stoeger, T. and Pelkmans, L.** (2013). Image-based transcriptomics in thousands of single human cells at single-molecule resolution. *Nat. Methods* **10**, 1127-1133.
- Bennett, J. T., Stickney, H. L., Choi, W.-Y., Ciruna, B., Talbot, W. S. and Schier, A. F.** (2007). Maternal nodal and zebrafish embryogenesis. *Nature* **450**, E1-E2.
- Besse, F. and Ephrussi, A.** (2008). Translational control of localized mRNAs: restricting protein synthesis in space and time. *Nat. Rev. Mol. Cell Biol.* **9**, 971-980.
- Boettiger, A. N. and Levine, M.** (2013). Rapid transcription fosters coordinate snail expression in the *Drosophila* embryo. *Cell Rep.* **3**, 8-15.
- Breiman, L.** (2001). Random forests. *Machine Learning* **45**, 5-32.
- Gross-Thebing, T., Paksa, A. and Raz, E.** (2014). Simultaneous high-resolution detection of multiple transcripts combined with localization of proteins in whole-mount embryos. *BMC Biol.* **12**, 55.
- Grün, D., Kester, L. and van Oudenaarden, A.** (2014). Validation of noise models for single-cell transcriptomics. *Nat. Methods* **11**, 637-640.
- Harvey, S. A. and Smith, J. C.** (2009). Visualisation and quantification of morphogen gradient formation in the zebrafish. *PLoS Biol.* **7**, e1000101.
- Harvey, S. A., Tumpel, S., Dubrulle, J., Schier, A. F. and Smith, J. C.** (2010). no tail integrates two modes of mesoderm induction. *Development* **137**, 1127-1135.
- Hoyle, N. P. and Ish-Horowicz, D.** (2013). Transcript processing and export kinetics are rate-limiting steps in expressing vertebrate segmentation clock genes. *Proc. Natl. Acad. Sci. USA* **110**, E4316-E4324.
- Itzkovitz, S., Lyubimova, A., Blat, I. C., Maynard, M., van Es, J., Lees, J., Jacks, T., Clevers, H. and van Oudenaarden, A.** (2011). Single-molecule transcript counting of stem-cell markers in the mouse intestine. *Nat. Cell Biol.* **14**, 106-114.
- Itzkovitz, S., Blat, I. C., Jacks, T., Clevers, H. and van Oudenaarden, A.** (2012). Optimality in the development of intestinal crypts. *Cell* **148**, 608-619.
- Jambor, H., Surendranath, V., Kalinka, A. T., Meijstrik, P., Saalfeld, S. and Tomancak, P.** (2015). Systematic imaging reveals features and changing localization of mRNAs in *Drosophila* development. *Elife* **4**, e05003.
- Junker, J. P., Noël, E. S., Guryev, V., Peterson, K. A., Shah, G., Huisken, J., McMahon, A. P., Berezikov, E., Bakkers, J. and van Oudenaarden, A.** (2014a). Genome-wide RNA Tomography in the zebrafish embryo. *Cell* **159**, 662-675.
- Junker, J. P., Peterson, K. A., Nishi, Y., Mao, J., McMahon, A. P. and van Oudenaarden, A.** (2014b). A predictive model of bifunctional transcription factor signaling during embryonic tissue patterning. *Dev. Cell* **31**, 448-460.
- Kim, D. H., Grün, D. and van Oudenaarden, A.** (2013). Dampening of expression oscillations by synchronous regulation of a microRNA and its target. *Nat. Genet.* **45**, 1337-1344.
- Kimmel, C. B., Warga, R. M. and Schilling, T. F.** (1990). Origin and organization of the zebrafish fate map. *Development* **108**, 581-594.
- Kimmel, C. B., Ballard, W. W., Kimmel, S. R., Ullmann, B. and Schilling, T. F.** (1995). Stages of embryonic development of the zebrafish. *Dev. Dyn.* **203**, 253-310.
- Lécuyer, E., Yoshida, H., Parthasarathy, N., Alm, C., Babak, T., Cerovina, T., Hughes, T. R., Tomancak, P. and Krause, H. M.** (2007). Global analysis of mRNA localization reveals a prominent role in organizing cellular architecture and function. *Cell* **131**, 174-187.
- Levesque, M. J. and Raj, A.** (2013). Single-chromosome transcriptional profiling reveals chromosomal gene expression regulation. *Nat. Methods* **10**, 246-248.
- Little, S. C., Tikhonov, M. and Gregor, T.** (2013). Precise developmental gene expression arises from globally stochastic transcriptional activity. *Cell* **154**, 789-800.
- Lyubimova, A., Itzkovitz, S., Junker, J. P., Fan, Z. P., Wu, X. and van Oudenaarden, A.** (2013). Single-molecule mRNA detection and counting in mammalian tissue. *Nat. Protoc.* **8**, 1743-1758.
- Mueller, F., Senecal, A., Tantale, K., Marie-Nelly, H., Ly, N., Collin, O., Basyuk, E., Bertrand, E., Darzacq, X. and Zimmer, C.** (2013). FISH-quant: automatic counting of transcripts in 3D FISH images. *Nat. Methods* **10**, 277-278.
- Nair, G., Walton, T., Murray, J. I. and Raj, A.** (2013). Gene transcription is coordinated with, but not dependent on, cell divisions during *C. elegans* embryonic fate specification. *Development* **140**, 3385-3394.
- Oka, Y. and Sato, T. N.** (2015). Whole-mount single molecule FISH method for zebrafish embryo. *Sci. Rep.* **5**, 8571.
- Ollion, J., Cochenne, J., Loll, F., Escude, C. and Boudier, T.** (2013). TANGO: a generic tool for high-throughput 3D image analysis for studying nuclear organization. *Bioinformatics* **29**, 1840-1841.
- Padovan-Merhar, O., Nair, G. P., Biaisch, A. G., Mayer, A., Scarfone, S., Foley, S. W., Wu, A. R., Churchman, L. S., Singh, A. and Raj, A.** (2015). Single mammalian cells compensate for differences in cellular volume and DNA copy number through independent global transcriptional mechanisms. *Mol. Cell* **58**, 339-352.
- Pauli, A., Valen, E., Lin, M. F., Garber, M., Vastenhouw, N. L., Levin, J. Z., Fan, L., Sandelin, A., Rinn, J. L., Regev, A. et al.** (2012). Systematic identification of long noncoding RNAs expressed during zebrafish embryogenesis. *Genome Res.* **22**, 577-591.
- Peterson, K. A., Nishi, Y., Ma, W., Vedenko, A., Shokri, L., Zhang, X., McFarlane, M., Baizabal, J.-M., Junker, J. P., van Oudenaarden, A. et al.** (2012). Neural-specific Sox2 input and differential Gli-binding affinity provide context and positional information in Shh-directed neural patterning. *Genes Dev.* **26**, 2802-2816.
- Preibisch, S., Saalfeld, S. and Tomancak, P.** (2009). Globally optimal stitching of tiled 3D microscopic image acquisitions. *Bioinformatics* **25**, 1463-1465.
- Raj, A., van den Bogaard, P., Rifkin, S. A., van Oudenaarden, A. and Tyagi, S.** (2008). Imaging individual mRNA molecules using multiple singly labeled probes. *Nat. Methods* **5**, 877-879.
- Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. and Regev, A.** (2015). Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**, 495-502.
- Schier, A. F. and Talbot, W. S.** (2005). Molecular genetics of axis formation in zebrafish. *Annu. Rev. Genet.* **39**, 561-613.
- Schindelin, J., Arganda-Carreras, I., Frise, E., Kaynig, V., Longair, M., Pietzsch, T., Preibisch, S., Rueden, C., Saalfeld, S., Schmid, B. et al.** (2012). Fiji: an open-source platform for biological-image analysis. *Nat. Methods* **9**, 676-682.
- Thisse, C. and Thisse, B.** (2008). High-resolution in situ hybridization to whole-mount zebrafish embryos. *Nat. Protoc.* **3**, 59-69.
- Tomancak, P., Berman, B. P., Beaton, A., Weiszmänn, R., Kwan, E., Hartenstein, V., Celniker, S. E. and Rubin, G. M.** (2007). Global analysis of patterns of gene expression during *Drosophila* embryogenesis. *Genome Biol.* **8**, R145.
- Tu, Z. and Bai, X.** (2010). Auto-context and its application to high-level vision tasks and 3D brain image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **32**, 1744-1757.
- Vargas, D. Y., Raj, A., Marras, S. A. E., Kramer, F. R. and Tyagi, S.** (2005). Mechanism of mRNA transport in the nucleus. *Proc. Natl. Acad. Sci. USA* **102**, 17008-17013.
- Weiser, D. C., Pyati, U. J. and Kimelman, D.** (2007). Gravin regulates mesodermal cell behavior changes required for axis elongation during zebrafish gastrulation. *Genes Dev.* **21**, 1559-1571.

Supplementary Materials and Methods

Below we describe the image analysis tools that we developed and made available in Fiji and KNIME. A visual representation of the flow of the image analysis pipeline can be found in Fig. S6. Here, we follow the order in which the tools are described in the paper.

Transcript detection

To analyze transcripts detected by smFISH we developed a Fiji plugin, which we called *Transcript analysis*. Transcripts are sub-resolution structures that appear in 3D image stacks as sharp peaks of signal over background. To remove background, images were filtered using tophat filtering. The filter radius was chosen equal to the maximum transcript radius, R_{tx} , which equals $0.3\ \mu\text{m}$ in our images, such that larger structures were filtered out. Next, the images were smoothed with a Gaussian kernel of standard deviation σ_{tx} to remove the noise. The value of σ_{tx} was fixed to $0.1\ \mu\text{m}$ (1 pixel) to avoid the transcript signal being smoothed out. To distinguish transcripts from background signal, an intensity threshold T_{tx} was determined. To choose T_{tx} , the intensity distribution of maxima that remained after filtering was plotted. One or two very sharp peaks corresponding to the background were visible at low intensities. At higher intensities, the maxima corresponding to transcripts formed an additional lower peak. T_{tx} was chosen between the background peak(s) and the transcripts peak to detect transcripts and exclude noise. To prevent over-segmentation of transcripts, a minimum distance equal to R_{tx} , the maximum transcript radius, was imposed between local maxima. Next, transcripts were segmented using watershed segmentation initiating from the detected maxima. For the watershed we used the mcib package implementation that can be installed within Fiji through the *ImageJ 3D suite* update site (Ollion et al., 2013). The watershed process was stopped at $2/3$ of T_{tx} (stopping at T_{tx} would lead to very small segmentations for the detected transcripts with intensity values close to T_{tx}). To avoid over-segmentation of large spots (e.g. transcription foci) an additional round of filtering and detection was applied. For this, single transcripts were removed through filtering, after which a threshold of value T_{tx} , the transcript detection threshold, was applied again to the remaining signal. The regions defined after this thresholding were used to fuse the signal that was incorrectly split during transcript segmentation. Finally, a 3D labeled image was obtained where each transcript region was identified by a unique integer value. These regions include single transcripts and transcription foci.

The Transcript analysis plugin is available through the Fiji update site MS-ECS-2D.

Detection of transcription foci

Transcript detection results in the segmentation of both transcripts and transcription foci without distinguishing between the two. Transcription foci differ from transcripts in size and sum intensity. To detect transcription foci, the maximum intensity and volume of each segmented region was calculated and regions with an intensity and volume larger than thresholds T_{imax} and T_{vol} , respectively, were selected as foci candidates. Each threshold was defined as the median of the feature plus 3 times the median absolute deviation (MAD) of the feature. Thresholds were tuned for each probe set by multiplying the MAD with a value between 1 and 3, depending on the intensity and volume characteristics for transcription foci of that gene. To exclude rare background spots that result from non-dissolved probe clusters, a focus candidate was only considered a true focus if its average position overlapped with a nucleus. To allow the detection of foci that were located at the edge of a nucleus we dilated the nuclear detection with a set pixel value of 10. Finally the number of transcripts in a transcription focus was calculated as the ratio of focus integrated intensity to median integrated intensity of all transcripts in the image as described previously (Mueller et al., 2013).

The detection and quantitative analysis of transcription foci is part of the Transcript analysis plugin.

Nucleus segmentation

Nuclei were segmented based on their staining with DAPI using a watershed-based approach (Ollion et al., 2013). Segmentation was done in 2D on a maximum projection of the original image along the z-axis of the image stack. First, the value R_{nuc} was used to filter out background signal. This value was chosen to represent the maximum nuclei radius in our images, typically 6 μm . Next, nuclei were detected and segmented using watershed segmentation. The stopping threshold for the watershed segmentation was determined automatically based on the Otsu threshold.

The segmentation of nuclei is part of the Transcript analysis plugin.

Cell segmentation

Overview

Segmentation of cells was based on Phalloidin staining and consists of two steps. First a pixel-level classifier was used to assign to each pixel a probability of being on a membrane, as well as a probability of being at the intersection of multiple membranes (i.e. a *vertex*). To this end we trained a cascaded Random Forest classifier. Our membrane and vertex classification is available in KNIME. Based on the output of this classifier, we traced membranes as highly likely paths between vertices. For this we developed a plugin in Fiji, which we called *PathFinder*. In the following we give a short overview on Random Forest classifiers as well as the concept of cascading. Then we describe the specific settings of our cascaded Random Forest implementation as used to classify membranes and vertices, followed by a description of our membrane-tracing algorithm *PathFinder*.

Cascaded Random Forests for Pixel Classification

Random Forest classifiers (Breiman, 2001) are vastly applied in computer vision as well as biomedical image analysis. The basic idea behind Random Forests is to train a set of decision trees from manually segmented training data. A decision tree captures a hierarchy of threshold-based yes/no decisions, one per node of the tree. Thresholds are applied to feature images selected from a large bank of image filters, like e.g. Gauss filters, image derivatives and Laplacians, and Gabor filters. If a pixel passes a threshold, it is passed down the left branch of the tree below the respective node, and if it falls below the threshold, it is passed down the right branch. Thresholds as well as feature selection are learned automatically from training data subject to the objective of maximizing an information theoretic measure of *information gain* in each node. The leaves of a tree encode class probabilities of pixels as relative fractions of training pixels of a given class that are passed down to the respective leaf. A forest is a collection of trees. A forest predicts class probabilities by averaging the probabilities estimated by individual trees.

Cascaded Random Forests (Tu and Bai, 2010), also referred to as Auto Context models, have recently become popular in computer vision. The idea is to apply a cascade of Random Forests sequentially. The forest at the top level of the cascade is a classical Random Forest, as described above. Forests at subsequent layers of the cascade receive additional input feature images, namely

the probability images yielded by the respective previous layer. This approach has been shown to produce considerably smoother and more accurate segmentations than single Random Forests.

Membrane and vertex classification

We adopted the cascaded Random Forest approach for pixel-wise classification of images into membrane, vertex, and background. As for the parameters of our cascade of Random Forests: We trained a two-level cascade. As bank of image filters, we used the filters provided by the Fiji *Trainable Weka Segmentation* plugin.

In addition, we used the Watershed segmentation as a feature image, as it is informative in regions with uninterrupted membrane signal. Furthermore, we used difference feature images. These contain differences between feature images and translated copies thereof. We set the maximum translation to about one fifth of an average membrane length, which equals 30 pixels in our images. Difference feature images allow for learning informative contextual features not yet captured by the pre-defined filter bank. Each forest was composed of 16 decision trees of depth 12. Above depth 12, a node was declared a leaf during training if it received less than 5 training pixels. Training pixels were split into left and right branches such as to optimize the *gini information gain*. We trained our cascade of forests on a set of distinct training images for which we generated “ground truth” segmentations manually. Given an input image, the trained cascade of forests generated two output images: A membrane probability map, and a vertex probability map.

The membrane and vertex classification plugin is available in KNIME (<http://tinyurl.com/KNIME-MS-ECS>).

Membrane tracing

Given membrane and vertex probability maps, the *PathFinder* plugin proceeds as follows: (1) Each pixel whose probability of being a vertex exceeds a threshold t_v is classified as a vertex pixel. In general, vertex pixels form connected components. The centers of mass of these connected components serve as vertex locations. (2) A distance transform of the image is computed as follows: Each pixel is assigned a *membrane cost*, namely its probability of *not* being membrane; for each pixel, the path with cheapest sum of costs to a vertex location is computed, and this sum of costs is

stored at the pixel. (3) The watershed segmentation of the distance transform image is computed, where we only consider pixels with a membrane probability above a threshold t_m . (4) Pairs of vertices whose watershed regions touch are connected by respective cheapest paths. (5) The average membrane cost of each path is computed, namely the sum of pixel-wise membrane cost divided by the path length. Paths with an average membrane cost above t_p are removed. Cells are identified as the compartments of the image tessellation defined by the remaining paths. (6) Small cells with weak membrane probabilities are removed as follows: For each cell, the most expensive adjacent path is identified. Let a denote the area of the cell, and c the average membrane cost of the most expensive adjacent path, then the path is removed if $c > \log(a+1) * t_a$, where t_a is a parameter of our method. The remaining paths constitute the automated membrane segmentation. We set the parameters of our method as follows: $t_v = 40\%$, $t_m = 1,5\%$, $t_p = 50\%$, and $t_a = 0.053$. The final output of the plugin is a cell segmentation mask.

We performed a two-fold quantitative evaluation of our method in terms of precision (fraction of obtained cell segmentations that is correct; lower values indicate more over segmentation) and recall (fraction of total cells that is correctly segmented; lower values indicate more under segmentation) on 11 images of zebrafish animal caps (Fig. S6). We studied how varying the parameter t_a affects the performance of our method and the resulting ROC curve is shown in Fig. S6A. Our choice of parameters favors precision over recall, because we found manual corrections to be more efficient in this case, i.e. missed membranes are immediately visible to the biologist when the segmentation result is overlaid onto the phalloidin staining image, whereas false positive membrane detections can only be found by switching this overlay off and on. At the chosen parameters our method achieves an average precision of 91% and recall of 70% (Fig. 6B).

The PathFinder plugin (to trace membranes) is available through the Fiji update site MS-ECS-2D.

Our membrane segmentation pipeline has some similarities with a previously developed approach (Cilla et al., 2015), as both rely on detecting vertices and tracing membranes between them. However, in contrast to Cilla et al., we use cascaded Random Forests instead of single image filter to obtain vertex and membrane probabilities. Thus, we automatically learn from training data which features are best for discriminating membrane, vertices, and background.

Optimizing cell segmentation and annotating cells

To facilitate correction of the cell segmentation output of the *PathFinder* plugin, we developed a tool in Fiji called *Cell annotation*. The tool allows splitting and fusing of cells by simply drawing or erasing segmentation contours in one click, using a maximum projection of the membrane channel as a reference. The image region is automatically relabeled according to the updated cell boundaries. In a second mode, cells can be assigned to different cell types by clicking on them. In our images we distinguished between DEL, EVL and YSL as well as regions outside of the embryo. The output of the *Cell annotation* tool is a dual channel image with the corrected cell segmentation in the first channel and the cell type mask -if applicable- in the second channel.

The Cell annotation plugin is available through the Fiji update site MS-ECS-2D.

Integration of all components of smFISH image analysis (also see Fig. S6)

Cell and tissue masks can be used as input for the *Transcript analysis* plugin to enable the analysis of transcript abundance at cellular resolution. In the *Transcript analysis* plugin, nuclei are automatically assigned to cell regions with the restriction that each cell can contain maximally one nucleus. The detected transcripts and transcription foci are assigned to cells and, if applicable, to nuclei. This information is used to calculate the number of transcripts per cell and per nucleus. For each cell and each nucleus the area is determined to calculate transcript density.

Because nuclei were segmented in 2D, a small nuclear size in the z-projection may indicate incomplete presence of the nucleus in the z-stack. This might result in the assignment of cytoplasmic transcripts to the nucleus, and conversely, to the dilution of nuclear transcript density by the absence of transcripts in the cytoplasm. In our experiments, nearly 70% of nuclei are present in all slices of the imaged z-stack, 40% of which is centered in the z-slice (Fig. S11). To prevent errors in transcript density calculations due to nuclei that are present only in a small portion of the z-stack, the surface area of nuclei in maximum projections was used as a proxy for their representation in the z-stack, and only those cells with the top 25% largest nuclei were analyzed, because they are most likely to be centered in the z-section. Finally, because in these cells the nucleus makes up a

larger proportion of the total number of pixels than in the whole animal cap, we corrected the obtained transcript counts in each class for the nuclear representation in the whole animal cap. When calculating subcellular transcript distribution in other systems, the analyzed sample thickness can be adjusted to the size of the nuclei to ensure that a large enough proportion of nuclei is positioned centrally in the z-stack.

Supplementary References

Breiman, L. (2001). Random forests. *Machine Learning* **45**, 5–32.

Cilla, R., Mechery, V., Hernandez de Madrid, B., Del Signore, S., Dotu, I. and Hatini, V. (2015). Segmentation and tracking of adherens junctions in 3D for the analysis of epithelial tissue morphogenesis. *PLoS Comput Biol* **11**, e1004124.

Mueller, F., Senecal, A., Tantale, K., Marie-Nelly, H., Ly, N., Collin, O., Basyuk, E., Bertrand, E., Darzacq, X. and Zimmer, C. (2013). FISH-quant: automatic counting of transcripts in 3D FISH images. *Nat Meth* **10**, 277–278.

Ollion, J., Cochenec, J., Loll, F., Escude, C. and Boudier, T. (2013). TANGO: a generic tool for high-throughput 3D image analysis for studying nuclear organization. *Bioinformatics* **29**, 1840–1841.

Pauli, A., Valen, E., Lin, M. F., Garber, M., Vastenhouw, N. L., Levin, J. Z., Fan, L., Sandelin, A., Rinn, J. L., Regev, A., et al. (2012). Systematic identification of long noncoding RNAs expressed during zebrafish embryogenesis. *Genome Research* **22**, 577–591.

Tu, Z. and Bai, X. (2010). Auto-context and its application to high-level vision tasks and 3D brain image segmentation. *IEEE Trans Pattern Anal Mach Intell* **32**, 1744–1757.

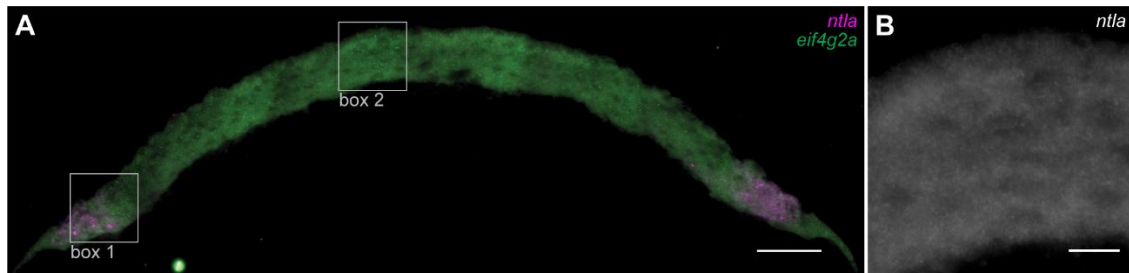


Fig. S1. smFISH for *ntl* and *eif4g2a* at 50% epiboly. (A) smFISH for *ntl* and *eif4g2a* on a cross-section of an embryo at 50% epiboly showing the complete view of the animal cap shown in Fig. 1A. *ntl* is expressed in the margin. *eif4g2a* is ubiquitously expressed. Box 1: region depicted in Figs 1B-D. Box 2: region depicted in Fig. S2B. Scale bar: 50 μ m. (B) Detail of (A) showing the absence of *ntl* transcripts at the animal pole. Scale bar: 10 μ m. Images are maximum projections of 17 z-slices spaced by 0.3 μ m.

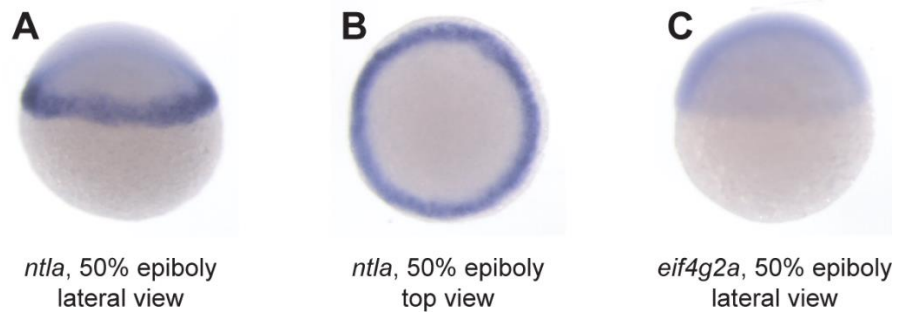


Fig. S2. Whole mount *in situ* hybridization for *ntlA* and *eif4g2a* at 50% epiboly. (A,B) Whole mount *in situ* hybridization for *ntlA* at 50% epiboly stage, lateral view (A) and top view (B). *ntlA* is expressed in the margin of the embryo. (C) Whole mount *in situ* hybridization for *eif4g2a* at 50% epiboly stage, side view. *eif4g2a* is ubiquitously expressed.

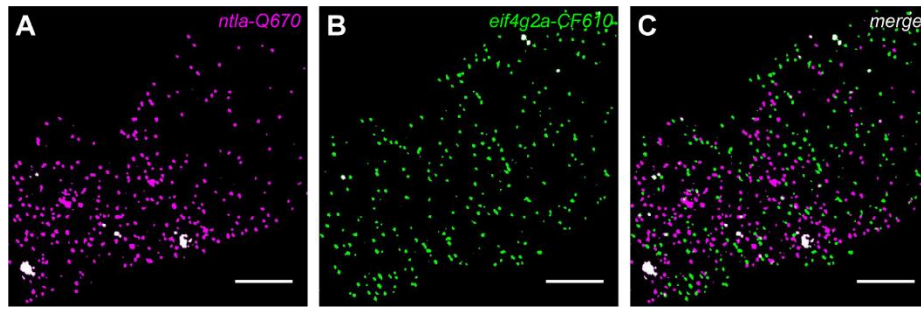


Fig. S3. Dual color detection of *ntlA* and *eif4g2a* by smFISH does not display overlap. Dual color detection of transcripts for *ntlA* (A) and *eif4g2a* (B) in a single z-slice of Fig. 1B. (C) Merged image of the detections in (A) and (B). Transcript detection of the two probe sets does not overlap, indicating the specificity of our smFISH protocol. Scale bars: 10 μ m.

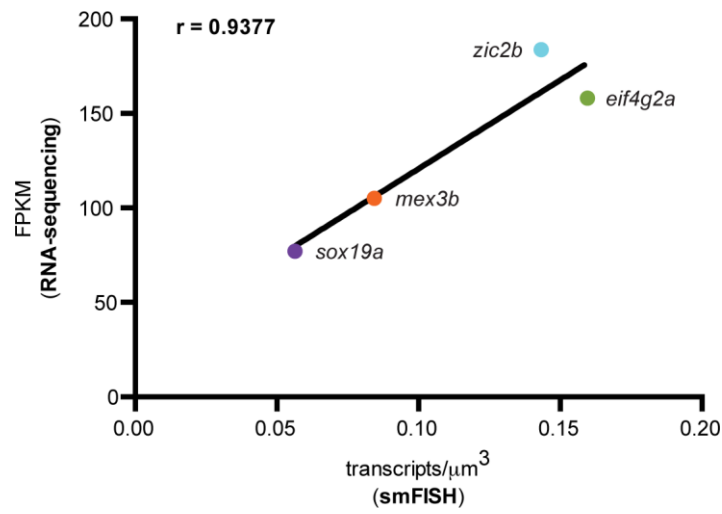


Fig. S4. Quantification of transcript numbers by smFISH correlates well with RNA-seq data. Correlation between average transcript density per embryo as measured by smFISH and relative expression in FPKM as measured by RNA-seq (Pauli et al., 2012) for four genes at dome stage. The correlation coefficient $r = 0.94$. smFISH data is an average of 3 embryos.

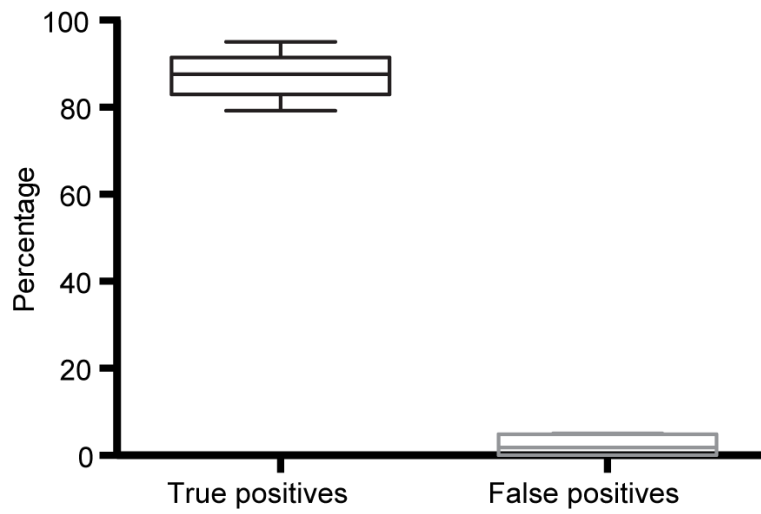


Fig. S5. Performance evaluation of automated transcription foci detection. Evaluation of transcription focus detection efficiency by comparing manual detection with automated detection. Six different images with a total of 256 transcription foci were analyzed. Shown is the percentage of true and false positives.

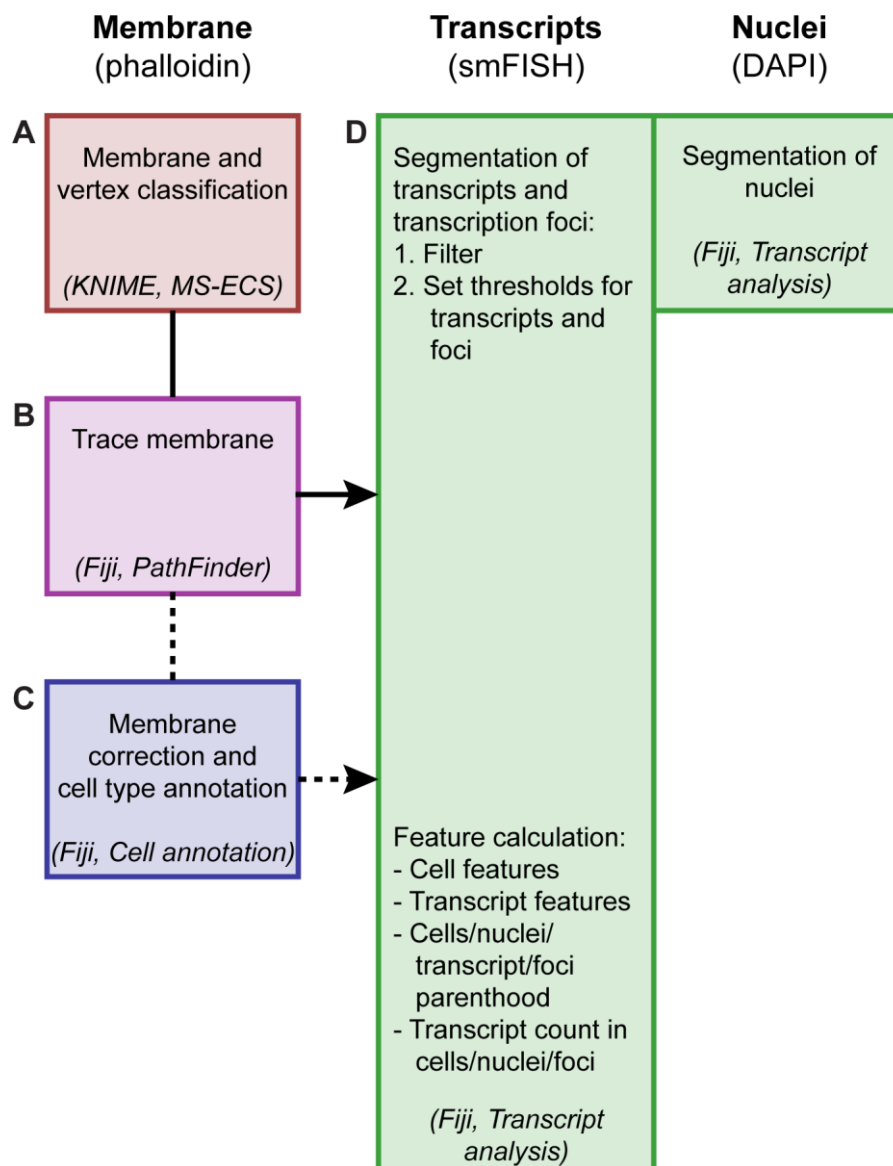


Fig. S6. Flow chart of the smFISH image analysis pipeline. (A) Cells are segmented based on a phalloidin-stained membrane image. Each pixel is assigned a probability of being part of the membrane, as well as a probability of being part of an intersection of multiple membranes (i.e. a vertex) with the *MS-ECS-2D KNIME* pipeline (B) Next, membranes are traced on the membrane-classified image as highly likely paths between vertices with the *PathFinder* plugin. The result of the cell segmentation is a 2D cell mask. (C) Optionally, the cell mask can be corrected with the *Cell annotation* plugin. This plugin can also be used to annotate cell types within the cell segmentation mask. The plugin outputs a (corrected) 2D cell segmentation and cell type mask as a multichannel image. (D) Nuclei and transcripts are segmented with the *Transcript analysis* plugin. The input for the nuclear segmentation is a 2D maximum projection of a z-stack image. A projection can be used because the samples are less than a cell layer thick, making it unlikely that two nuclei overlap in the image. Transcripts and transcription foci are segmented in 3D. After subtraction of the background and smoothing of the image with a Gaussian filter all maxima in the image are identified. The intensity distribution of these maxima is plotted to determine an appropriate threshold for transcript detection. Transcription foci are classified as foci based on their size and intensity. The segmentation results in a 3D mask of transcripts and foci. In the *Transcript analysis* pipeline, the cell labeling, nuclear labeling, transcript and transcription foci labeling, and optionally the cell type annotation are integrated to calculate a set of features. The final output of the analysis pipeline is a table of cell and transcript features that can be used to calculate e.g. cellular transcript density.

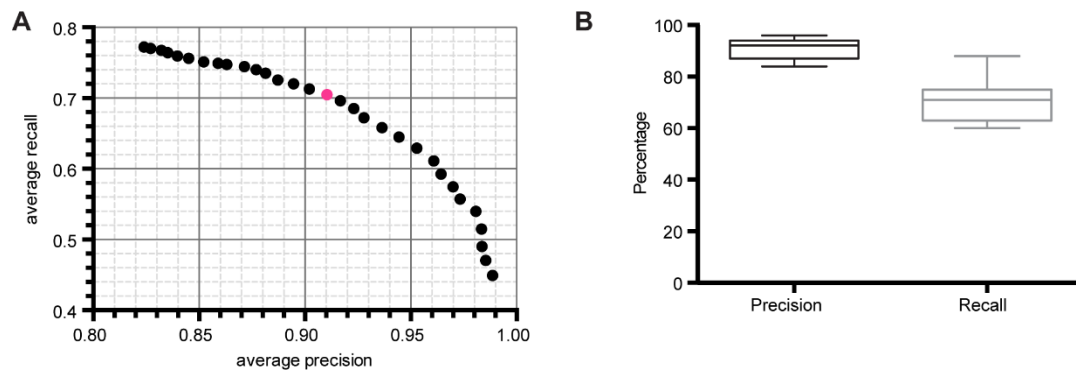


Fig. S7. Performance evaluation of the membrane segmentation pipeline. (A) Precision-recall curve obtained by varying the small cell removal parameter from 0.076 (top left) to 0.031 (bottom right). An increase in precision corresponds to a decrease in over segmentation. An increase in recall corresponds to a decrease in under segmentation. The plot shows averages over 11 images of zebrafish animal caps. The magenta dot corresponds to a parameter value of 0.053 for small cell removal, which results in the most efficient manual corrections. (B) Box plot displaying the spread of precision and recall values over 11 images for our chosen parameter set.

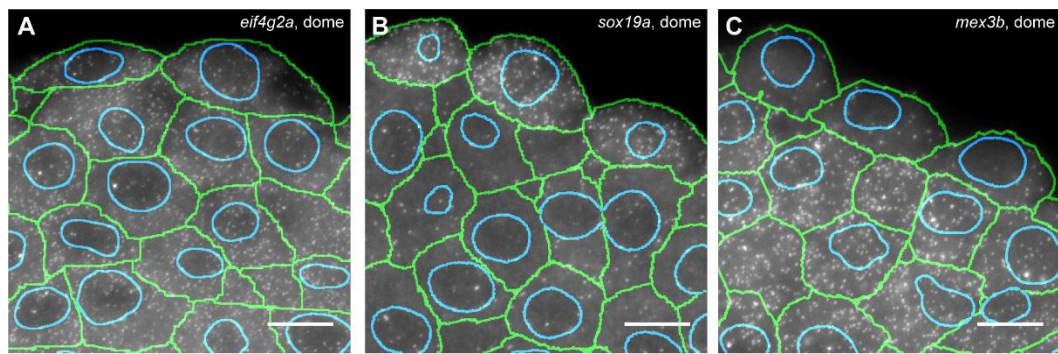


Fig. S8. Original smFISH images corresponding to Fig. 3. *sox19a* (A), *mex3b* (B) and *eif4g2a* (C) at dome stage. Green, cell outlines; blue, nuclear outlines. Scale bars: 10 μ m. smFISH detections are maximum projections of 17 z-slices spaced by 0.3 μ m.

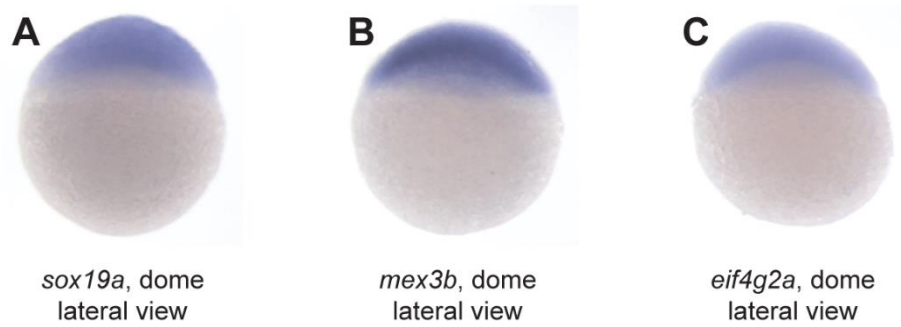


Fig. S9. Whole mount *in situ* hybridization for *sox19a*, *mex3b* and *eif4g2a* at dome stage. Whole mount *in situ* hybridization for *sox19a* (A), *mex3b* (B) and *eif4g2a* (C) at dome stage does not detect differences in expression levels between the enveloping layer (EVL) and the deep layer (DEL).

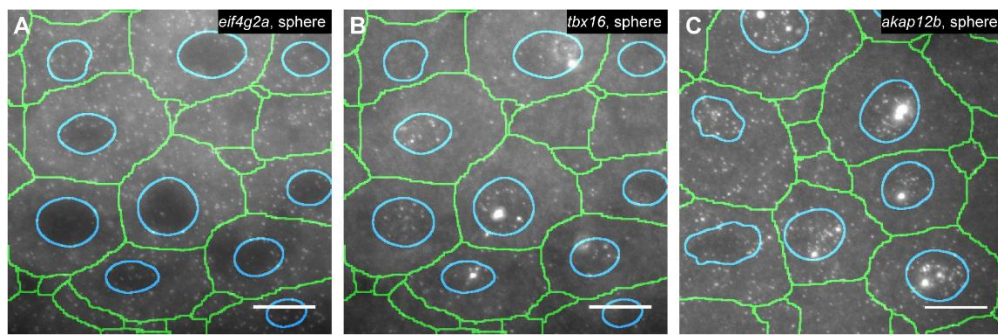


Fig. S10. Original smFISH images corresponding to Fig. 4. *eif4g2a* (A), *tbx16* (B), and *akap12b* (C) at sphere stage. Green, cell outlines; blue, nuclear outlines. Scale bars: 10 μm . smFISH detections are maximum projections of 17 z-slices spaced by 0.3 μm .

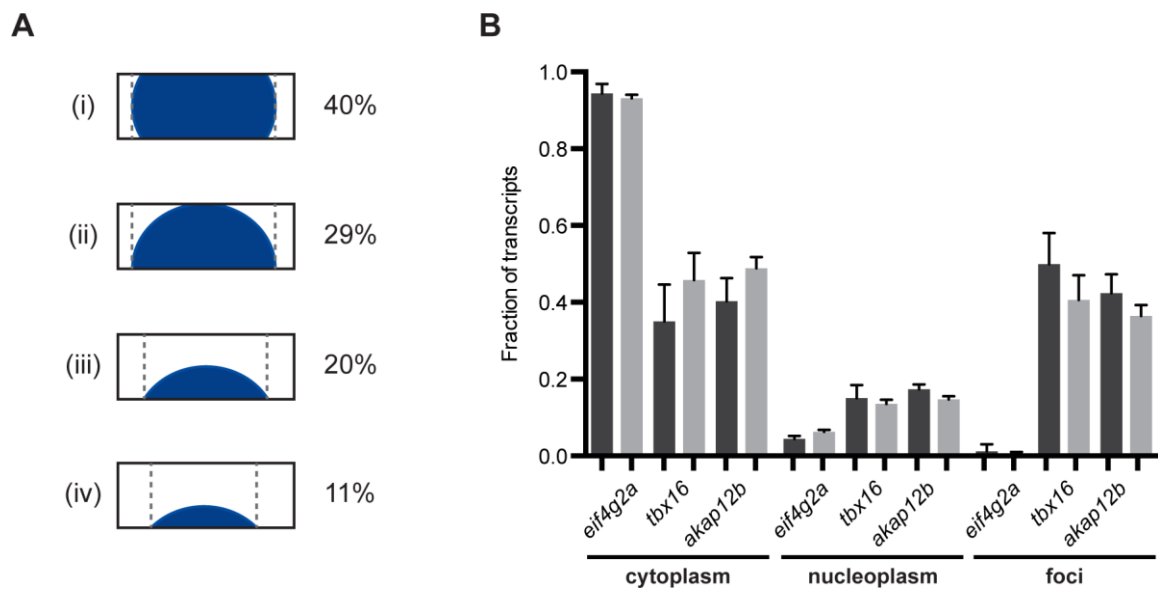


Fig. S11. Nuclear size selection improves quantification of subcellular transcript distribution.

(A) Nuclei were classified in four groups based on their presence in the z-stack. Group (i) represents nuclei that are present in all 17 slices of the z-stack, and that are centered in the middle (thus filling almost the complete segmented volume), (ii) represents nuclei that are present in all 17 slices of the z-stack, but are shifted to one side, (iii) represents nuclei that are present in at least half of the z-stack, but not in all z-slices, and (iv) represents nuclei that are present in less than half of the z-stack. The numbers on the right of the schematic represent the percentage of nuclei in the different groups and is based on our analysis of 330 nuclei from 3 embryos. Box, acquired z-stack; blue, nucleus; dashed line, nuclear segmentation in 2D. (B) Subcellular quantification of smFISH analysis in cells with the 25% largest nuclei in the maximum z-projection (dark gray), compared to all cells (light gray). Using the 25% largest nuclei improves results but does not change the general subcellular distribution of transcripts. Differences between the two approaches are small because when analyzing all cells, the largest segmented nuclei dominate the outcomes. Values represent averages of three embryos. Error bars represent s.e.m..

Table S1

<i>ntla-Quasar670</i>	<i>EIF4G-2a-CalFluor610</i>	<i>slc7a8a-CalFluor610</i>
aatcccggttgatactgttg	accgggatgaggaggagaac	gccaatgatgttaccgacaa
gagataagtccgacgatcct	cttgatgggttggttac	cgaactggcattttccaaca
ttgaggcagacatatttccg	atgagagtgaagggtgggg	caatcgctgtataatgcct
ctaaggagatgatccaggcg	gcactgaaacgagaagcacc	tgacataagagtagtgcct
tctgaaattcgctctccacg	aaggaacctatctgaacg	cagcaatccataatcgaga
gcgaaagttaatatcccgc	aagattgatcgtgtcgttc	agttggagaaggtgagagca
aaatttggccacaactccg	caaactttcaggggtcagt	aacaagcagatggcagcgag
tcatttcattggtgagctct	agggctttgtctacaatcag	gacatcctgtactctgtag
gaaacattcgctctccagtc	aaagtttgagcctcctctg	ccttgacagattgtacgatg
cattgcattagggctgagac	gaatgtagtgcctctgtttt	agcaatcagaccgacatcat
caaaaaccagcaggaccgag	aggtgagaggactgtcattt	attgaggaaattccagccac
gttcacgtattccaccgat	aatgagatcgagcttgccaa	agatgaagatagcacggggg
ctttgtgaaagatacgggt	ccatatacctaagttggact	gacgtaggcagatgttagcaa
cctccggttgagttattgga	tttctcatgatctagccttg	caacagcttctaccaaacg
ctgggtcgtattgtgcaa	aatcgaatcctagcaggcaa	caaatgtggacaaggccaca
accgactttcacgatgtgta	aaatccttactgacatcctg	cgagcaccgcgaaaaacaa
gactgctgatcattttctga	gtcggcatgaacgaattctc	tgaagagcagggtggaata
aactgtgtctcaggaaaaga	caaacatatccgctaattccc	tgagtgtgtacatgactg
ctgatatgctgtgactgcaa	ctgtgccattgaaagagg	gcaacagtgacaccgtagaa
tcagagcggtaatactctca	ttagactgagcttgcgttg	aattggtcgggtcatactg
cgaggaaagccttggcaaaa	gctaactctcatcagcattca	gctccgagtacagagagaag
gacttcttgtggtcacttc	ctttggcatttggttttgt	gccaagaatagacagggga
cagattgctggttgcagtg	cttctgatcgtctctgtg	ggatgtcatctgtcaaca
cagccaccgagttgtgaata	ctcaggtaggaagtgtttg	cagactcatcggtttcttt
ggagtatctctcacagtacg	gcttattgatcagagtgctg	tgattccagtgctgtttg
cagctctgtggttctcaag	acactgggtccagaacattga	tctctgttcagactctctc
cggaagagttgtccatgtag	acgagttcagcaatgactgc	agcagctgctaagaaaccag
ctgtcatgagacgcaagact	gaatctccaacatcgatct	ctagagcttttcttatctc
acagacttgggtactgactg	tcctttatccactttagat	ggtgttcttttaacccaa
taggaactgagatgacaggc	gcagataactgctcttcatc	atgtagtttggcagatctg
ccgagtaggacatcgaagaa	tccacatgatcatgcaggaa	ttggtgtggtggaatacact
ttaggcctggatgctacatt	cgcagtaacatttcttagg	cacgatagcaacttgcagct
tcagtagctctgagccacag	taatgtcttcttccatgct	attctgtctgaagtccacag
agtccctaaatgtgaagcga	aatggatgaaggcttggct	gcgtgaaagtgaactgtggg
atcaagtccataactgcagc	gcgtgaagtaatacgttggc	ctaactttgcaccactgttt
tagatttccctcctaagcca	ggtgggttaacagatgcta	tgatcatcccgaagctatta
accctgtttctgattgcaa	agcagcagaaaaggctcttg	tgatcaactgattctgtccc
atacgtatgcctcggatga	tgactgaggatcccgttctaa	tcattttggtgtgctgtttt
acagcttctatgcaatggtt	atgactccattacacacagc	actctgtgaccactacaatc
tgctgataacctgtaaccat	taattcatctctcccacatg	ctctccatagcttccaaa
aagggatagtttaacccaaa	taatccctgaacatgctgta	gcagtgaaattacaccgtgat
gccccaaaatgtatggctaa	gagacacatttctgtctgga	gtgcaagttttagactttt
agtgagtcaacagcacaatt	tacagacatcatgggcacac	tgaattggtaacgacaggcc
gatacaatgaaaccggacgt	caaatacatgtacagacgcc	caaagatgctcaatgcaca
ctcgttctacagaagcaca	ggccaaataaatgctggtgt	tttgttttagcacacggta
cgaaacagcaaagtctgtct	agttgaaaattgtcccaga	gtatcgccagtggtatgaa
cggtcacttttcaaagcgtg	cactatttagtgccaaccaa	agcttctctgtctgattt
aggacgaatagcagacaaaa	tacctacggaagcattgca	tccgtgtaagtggaacacg

<i>slc7a8a-Quasar670</i>
cctttgggctgacaaagatt
gaatccatacaatgagggcc
atlttggtatagtaacgcc
ctgcaagtctccaaagatg
ggtagtaggtagattacc
tgggaaatagaggctgcaga
ggaacaattcaccaagtca
aatgagtgccagaagttcc
atggttcgaaagcattggct
gaaggaccctgtagaatg
ttcacataagggctgaccag
cacgaaggtcacaagggaa
acagctacagcattagatgc
atgggcatgatccatgacat
gagaggaggtaacagggag
gaatcatagccagcagacta
gcacagcatcagaagtgtg
ggtagttgatgaaaccaca
ttaatccgcagcacaatct
ccagaagagcaggtagatca
cagtaagcatgatggccaag
aaacattggggctgtgtgc
accaccacacagaattctg
atgtgtgaaggagtgtgtct
acccctctttatacaaaag
aacagagatcgaggagcagc
ttgtaattctgtgacacgc
cctgtgctgtactttatac
tgtctcacagaatgtccaga
aaataaaacgagccctggcg
tctgaaccatcacaatcggt
acactgttacaacagccgag
ctgtaaagccggagacatga
accaaaccactgtgcaattt
aggctttcagatacatcaca
ccatccacacgattatcata
catctacgagtgaaggttt
acactgtatctgactgagca
caacacagaggaccaacat
aatgctttttatacagcgcc
aaatgtggatgggtgggagg
gaacccaaaaggcagtgttt
tggttattcagtcagtcgtg
ctcaggagcagcattacaac
gcacaatcagagtcgtatgt
accagaattgttttgctt
tcgcctcagattgtaagta
ctgatacaaatcgttcacc

<i>tbx16-Quasar670</i>
ttaagtccaatgctctggg
caaggctccatcagaactg
gatagcctgcattatttagc
gctgaaattgtgcttgaggt
cttgatggtaggaatcacgc
ggtttggatgatcatctc
aagctgatttgcagtgagg
agatgtacttgcatacggc
cttatactcagaccgtctt
tgagcttcagctgaggaag
tggaaacgtggatggtagcg
gtacagatcatcagcctgaa
aaaggttggaaagacgctcc
agtgcaggcagtaaaggagg
ctttggcaaagggtgtgg
ttgttccctcatctctaaa
gatcaggcagatttctgtt
atlttcttgcactcgtctc
cagagcttcacacgatgacg
ttcgtgaaatctgtacggct
gccaacgctggaagatgaag
tccaagactcgggactcaaa
ctgtaaagcggatgagttcc
gtgatagccatgatggacag
caccagcagcagatgagaa
agcatctcaaaatggctc
gttcacttttggtacactca
cttgcaagtgtagtttagtc
ctaaccctactatccaaata
aataatctgccagtggtgtt
ggaaaaccctgtgatttg
gggggtgtttgtgatagt
ttcaaatctgcactggcaga
aacgccactctgcaattatt
ccacaacactggcagattat
tcactggccaagattgttta
ctctcactaacgatttccga
gccaggattattttgcttat
aacattgagtgagttcggc
tatlttttcccttcacagt
catttccaaatggccatca
tttgggtgaaactaacgctt
accataaaaccataagca
atctttttgtgtccgttt
tgtgcattgacttccatttt
accctgtcattgtttatttg
tgtgggtccatttaagagg
tataactcagccacttcaa

<i>akap12b-Quasar670</i>
gagaatgaggagagcgggag
tggcactctattgtttc
tcttctgttttctctac
aaatcttctgaagcccact
cttctcattctgtctttct
ttcagcaggttttcagttt
tttcttgacatgttctca
cagttcagatggtttctct
cgtttcagtggtagttcag
agatcttgtgtctctcttt
gcttttggatgagagcttt
tgttcagcagtttgttcac
tctttcttttccaccatc
atltttcttttggctcttc
cagagcttacaagaggggac
ctttagctttctttctgtg
tcagctcttctactgttc
tggatgctgtgtaacacgtg
ctgttttctgttctttggg
ctcctgtgtttttctcaa
attgcattttctatgggctg
ttaatctctcaactgcgc
agtttctgctataggagtgt
acagattgcctctcaacag
caagttctgcactttctca
actgcaaattctcagtgcc
gctatatctgctgaatggga
agacaacctcaacctcattc
tctcagacacattttgctcg
actttgactctgctgatgt
gactattttagggggtttct
gcttggatttcatctgtgat
ttgttacctctgtttctga
acttctctgagatcagactc
tttctctggtacagctttt
tgtaacagatgtttctgccg
tacagctgtcagtggtactt
ccaacgftaaattcctgtgc
agtctcaggtacatcatctg
cagtgccatttgatgtcatc
atgacactgtgacaacctct
cagctctagtgcttcaat
tgatattggcagctgagtgat
accacacgtacatattctct
ttttcgtttgctacagctg
gctgagttttacctcatttt
tgaaccgctcatgtgtaca
tccttaatgctttcattggt

sox19a-Quasar670
cgactctttcagtggtgag
cctaaaagcgaagagggtcc
acaagttctgaagtgcggtt
ttgtaaacgcctgaacacaa
gccgtctgtaactaaaagca
aaatccactgacgctaaacg
taccggagatcacccgtaaa
ctgccaacagaagttagtgg
gacttcattgctgctccag
ctgtgttggtgctgtgacag
ttgtccataggatcgctact
accataaatgcggttcaggg
cattttggggttctcctgtg
cgcttgctgatttcagagtt
aataaagggctggttctccg
atctgggtactccttcatgt
gggtgtccttttcagcac
cggctgacagaggatactta
gagttggcaccattcatgta
gagtgctgctataggacatg
ttcaccatggacatgactg
gagaatgagacacctgctct
tcaggtctcctgttaaagga
tcctccagggatatacatgc
gttgagtagcctctttgtg
cctgatcagatgtgtgtgag
cccgaacaggaaggtaaag
tcaagttaatgttcacgcat
ttccagttcaaccatgtgt
ttctgaggcacttcaacaag
agaacaaagacatggacct
gggctcaatgacaaaatca
tacaagcagttccaacagtg
caaccaacggctagtttat
ttcaaccagtcaactgggtg
ctagctggttaccatgcaa
tgctttgaaatcagtgaga
accttatagcgttcacagc
aatgtgtgcttattccctt
tagatcatcaggacaagcag
agtcttatggtacgaaagct
atgaagagcagcacaatgtt
atgtctaaaccgtgtgcat
ggtgaacctcatggcataca
cagtcaaatcccctaagcg
ttctcagagcagttgagaa
aaccctactataagggcgt
gtcaaggtgtcaagtcaca

mex3b-Quasar670
gatattcccccttcgagtg
tggaataagcgttattccc
aatcagtgcaatcccgaat
aaacggctggtttaaaagc
actaggcatcttggataggt
cgttgctgacagggaaatt
attcagtcattgtcacgctc
cgacaatttcagcaaatgc
gcagggctttaatttgcaa
ccgcactggagtttgatat
cctgtcacaacaaagactgg
atcattgagaagtgtcagc
ctgttttatttcgagcgc
gacgatataaggtgtcgtct
tcaaacactggctctttgct
ggtctacattctcagggcatg
atctgtgaactcaatgagcc
gtaatcagccattctggagc
gttattgtgttgccgttgc
ccgctgtagacaaatccatt
gaaggtcaagtcagtgcaat
cgaggagtaggatcaaatc
attattggcgcgttcatact
attcgtagagaaaaccggtg
ccattattggcattcgtgga
acgacgatgatgaaggagc
attacttcgcttcgaagca
cgtatagcctgagtaactgc
gcacgttattaggtctgtt
tggatactgcagctctacat
cgctcccttttctgtgtg
atttctcccagctacagta
aacgctgtgtacagttgac
tgactgactaaggtacctcg
cgtgtatgctaacagtcgac
ggaaaaaaaaacgcaaggctc
tctgctaccagtacattcga
agttccatcgtgtgataagc
ggatagttgctgttctccc
taatacaggttctgtgccc
ccacatctcgctgatatta
ctgttcgatacactttgcgt
tctctcttttctgtaacgt
ggattcaaaagacgccttct
aacgggtccttttctgtga
tggtaaccactggaatctaa
taaacctatagtgccgcaac
tgaccatgaggttaaaggc

zic2b-CalFluor610
ttggatttagtagcgggtgc
catctgtacgattggtgagc
atgccacgtgatgaactgag
cttgcccgaatagggttaa
gtagctatatggctatgagc
actacacccgcgtttaaata
tcaccacgattttgaccaa
ccctaaaccagaaatctggt
cgccgattccataaaactca
aaatcgggggtggttcagttt
agaatagctgctgctgtg
gttaaaagtggagctggcga
gatgtggagctgaagatggg
tgggtgagagtgatggagag
attccggggaaaagtaggtg
tccattcaggatgtttggag
cgaaaacctcaccgggtaat
atattggctgtatggatcgg
ccatgttcatactcatgttt
ctgtttgatgcactgtttg
gggtcgatccacttacaat
aagttttgtgcagcatttt
atgtggttacactgttctgg
tggacattctccaaaagc
ttgctttgaaagttgtctc
gcactcgatagattcact
aaattctccgaacgcgcgaa
cacacaggaacggttttct
cagacgtgtgaacgtgcatg
agtttgacagatatggctt
tgtgtttctcaaggagctg
agatcagactggacggagtg
cagaatcaggcgacaagggtg
tggacgttaagctattgtgc
agtccttaaacgtaccactc
taggatatgcagttcttggc
ccctcattgtctatttaca
caccgacatgctgagaacat
ttctcggcttctttattca
aactgcctgcatttttagac
gcaacgcattagctacaaa
atacatccactctttgtctt
aacatcgtaacttgaccgga
ccccgaaaatactttcatgg
acatttgactgggcactta
ttgtgtccattttcactgta
tttctgcagattttccata
gtttggcatttttatcaca