

Quantitation in Colocalization Analysis: Beyond "Red+Green=Yellow"

Ming Yuan

Department of Statistics
Columbia University

ming.yuan@columbia.edu

<http://www.columbia.edu/~my2550>



(Joint work with Shulei Wang)

- ▶ High contrast
- ▶ High specificity (targeted molecules)
- ▶ High throughput
- ▶ Quantitative (fluorescence intensity, fluorescence lifetime etc.)
- ▶ High resolution ($\sim 20nm$)
- ▶ Dynamic (monitoring biological events for 24 hour)
- ▶

[*Zebrafish*; Cutrale, F. et al., 2017]

[*Drosophila*; Chhetri, R.K. et al., 2015]

What is Colocalization?

A Statistical View of Colocalization

Global Assessment of Colocalization

Local Identification of Colocalization

Concluding Remarks

What is Colocalization?

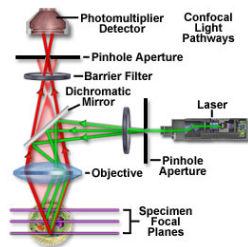
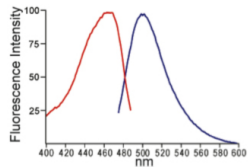
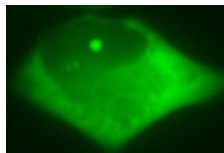
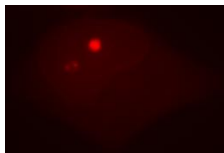
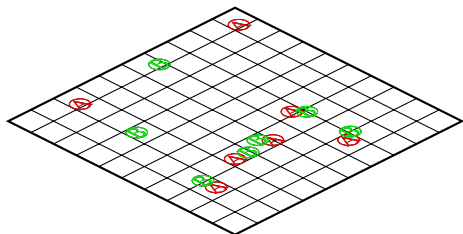
A Statistical View of Colocalization

Global Assessment of Colocalization

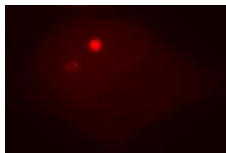
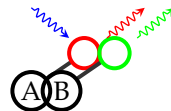
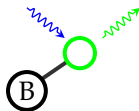
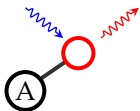
Local Identification of Colocalization

Concluding Remarks

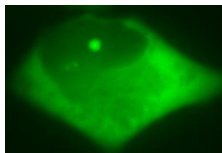
DUAL CHANNEL FLUORESCENCE IMAGING



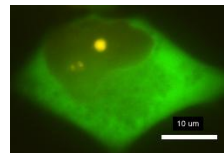
COLOCALIZATION



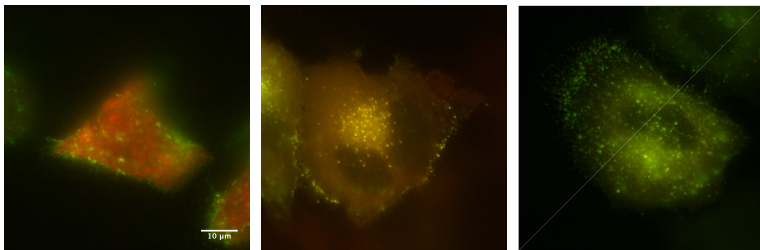
+



=

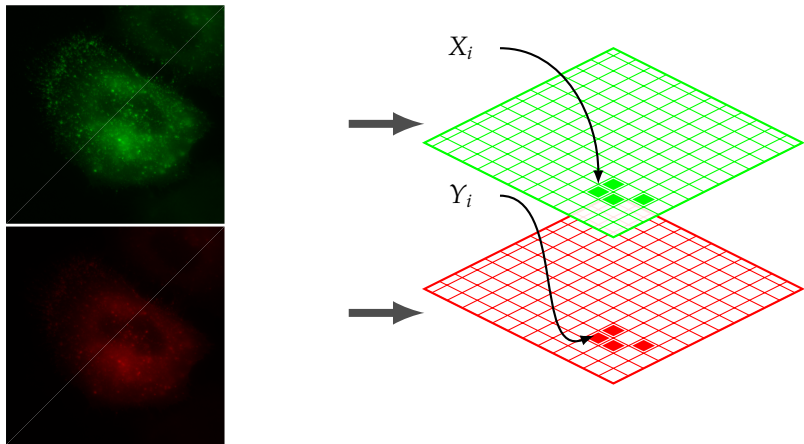


Red + Green = Yellow?



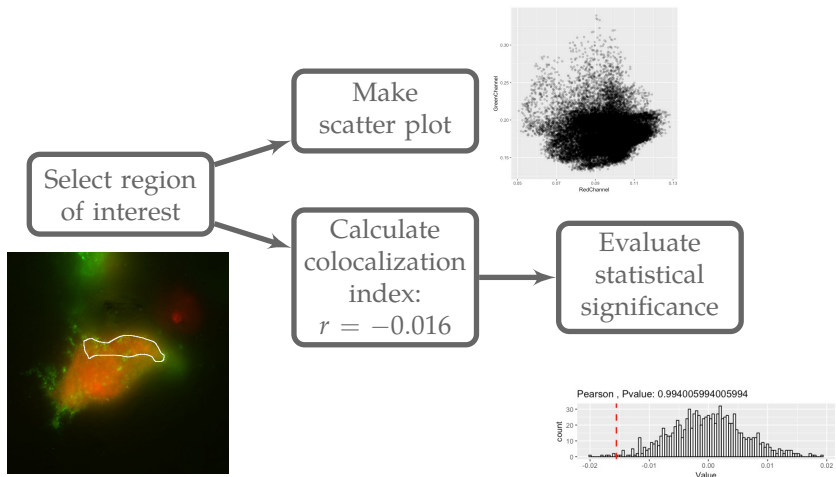
- ▶ Subjective, susceptible to cross-talk, and etc.
- ▶ Time consuming, labor intensive

PIXEL BASED MODELING



(Pioneered by Manders and Co., 1990s)

CURRENT PIPELINE FOR COLOCALIZATION



(see, e.g., Bolte and Cordelières, 2006, Dunn et al., 2011)

- ▶ How to choose region of interest?
- ▶ How to choose colocalization coefficient?
- ▶ How to evaluate statistical significance?
- ▶ How to do so in a computationally efficient way?

Our goal: a general *statistical/computational* framework for colocalization that is

- ▶ Automated
- ▶ Statistically valid
- ▶ Computationally efficient
- ▶ Flexible and powerful

What is Colocalization?

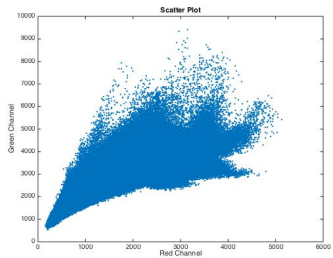
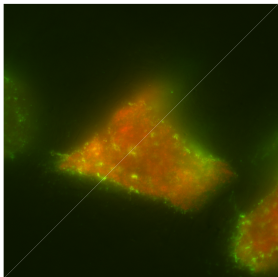
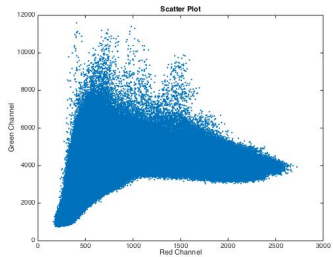
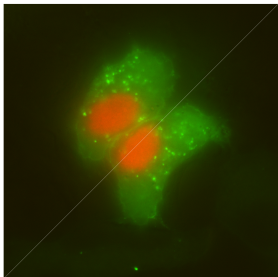
A Statistical View of Colocalization

Global Assessment of Colocalization

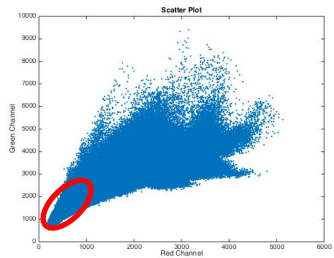
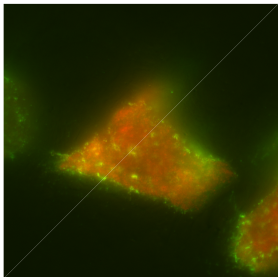
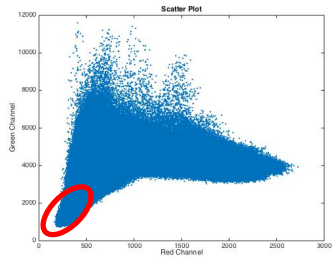
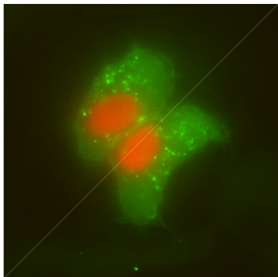
Local Identification of Colocalization

Concluding Remarks

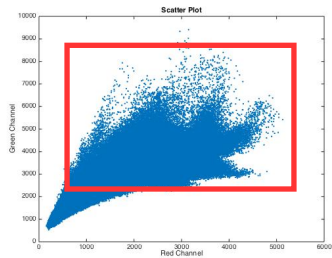
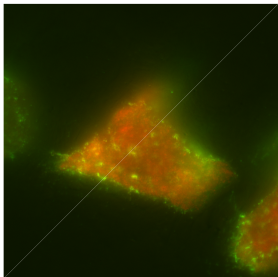
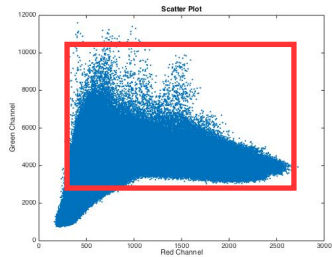
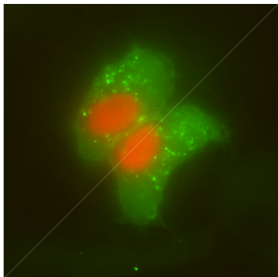
BACKGROUND OR SIGNAL?



BACKGROUND OR SIGNAL?



BACKGROUND OR SIGNAL?



WHAT ARE WE MEASURING?

Pearson's correlation coefficient:

$$r = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_i (X_i - \bar{X})^2 \sum_i (Y_i - \bar{Y})^2}}$$

Manders' colocalization coefficients:

$$M_1 = \frac{\sum_i X_i I_{(Y_i > 0)}}{\sum_i X_i}, M_2 = \frac{\sum_i Y_i I_{(X_i > 0)}}{\sum_i Y_i}$$

Correlation:



Co-occurrence:



- ▶ Positively quadrant dependence (PQD, for short) (Lehmann, 1966):

$$\mathbb{P}(X > x, Y > y) \geq \mathbb{P}(X > x)\mathbb{P}(Y > y)$$

- ▶ Colocalization manifested as correlated co-occurrent signals:

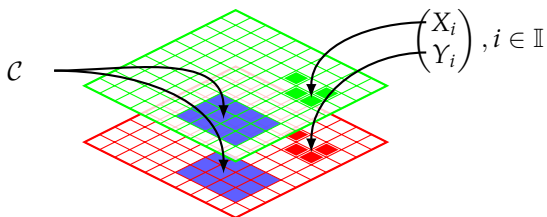
co-occurrence (V): $S(\eta_x, \eta_y) - S(\eta_x, -\infty)S(-\infty, \eta_y)$

correlation (T): $\mathbb{P}\{(X - \tilde{X})(Y - \tilde{Y}) > 0 | X, \tilde{X} > \eta_x; Y, \tilde{Y} > \eta_y\} -$
 $\mathbb{P}\{(X - \tilde{X})(Y - \tilde{Y}) < 0 | X, \tilde{X} > \eta_x; Y, \tilde{Y} > \eta_y\}.$

- ▶ Background vs signal:

$$F(x, y | x > \eta_x, y > \eta_y) = F_{\eta_x, \eta_y}(x, y) \quad \leftarrow PQD$$

COLOCALIZATION VIA STATISTICAL LENS



- ▶ Assume each $(X_i, Y_i)^\top$ is drawn from a bivariate distribution.
- ▶ Without colocalization

$$(X_i, Y_i) \sim \underbrace{F_0(x, y)}_{\text{no PQD}}$$

- ▶ With colocalization

$$(X_i, Y_i) \sim \underbrace{F_1(x, y)}_{\text{exhibit PQD}}$$

What is Colocalization?

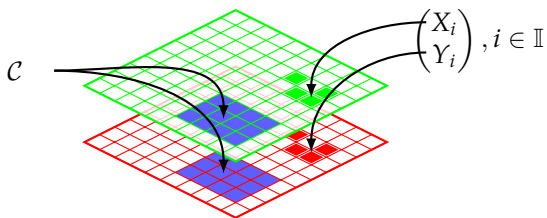
A Statistical View of Colocalization

Global Assessment of Colocalization

Local Identification of Colocalization

Concluding Remarks

A HYPOTHESIS TESTING APPROACH TO COLOCALIZATION



- ▶ Assume each $(X_i, Y_i)^T$ is drawn from a bivariate distribution.
- ▶ Without colocalization

$$(X_i, Y_i) \sim \underbrace{F_0(x, y)}_{\text{no PQD}}, \quad \forall i$$

- ▶ With colocalization located at an *unknown* set \mathcal{C} of pixels

$$(X_i, Y_i) \sim \underbrace{F_1(x, y)}_{\text{exhibit PQD}}, \quad \forall i \in \mathcal{C}$$

$$H_0 : (X_i, Y_i) \sim F_0 \quad \forall i \in \mathcal{C} \quad \text{vs} \quad H_1 : (X_i, Y_i) \sim F_1 \quad \forall i \in \mathcal{C}$$

- Positively quadrant dependent property implies

$$\tau_H := \mathbb{E}(\text{sign}(X_i - X_j)\text{sign}(Y_i - Y_j)) > 0.$$

Here τ_H is called *Kendall tau correlation*.

- Empirical version Kendall tau correlation is a good indicator of correlation of $H(x, y)$

$$\hat{\tau}_H := \frac{1}{n_{\mathcal{C}}(n_{\mathcal{C}} - 1)} \sum_{i \neq j \in \mathcal{C}} \text{sign}(X_i - X_j)\text{sign}(Y_i - Y_j)$$

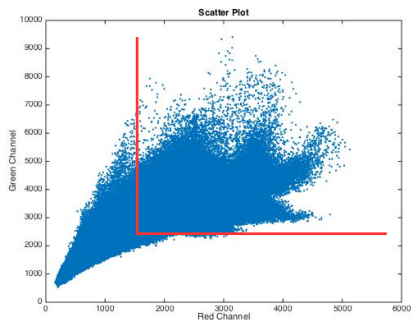
But we do *not* know \mathcal{C} (or equivalently η_x and η_y)!

TEST FOR CONDITIONAL PQD

- ▶ Known η_x and η_y – *conditioned and normalized* Kendall's tau

$$\hat{\tau}(\eta) = \begin{cases} \sqrt{\frac{18}{n_\eta(n_\eta-1)(2n_\eta+5)}} \sum_{i,j \in \mathcal{K}(\eta): i < j} \text{sign}(X_i - X_j)(Y_i - Y_j) & n_\eta > 1 \\ -\infty & n_\eta \leq 1 \end{cases}$$

where $\mathcal{K}(\eta) = \{i : X_i \geq \eta_x, Y_i \geq \eta_y\}$ and $n_\eta = |\mathcal{K}(\eta)|$.



- ▶ Unknown η_x and η_y ,

$$\tau^* := \max_{T_x \geq X_{(i)}, T_y \geq Y_{(i)} : i, j \geq \lfloor n/2 \rfloor} \hat{\tau}(T)$$

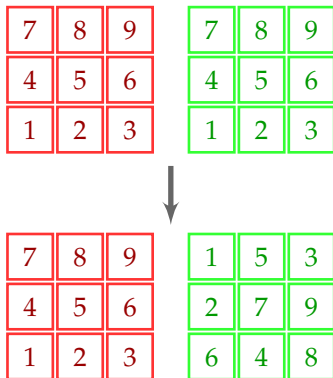
- ▶ Test

$$\psi_T = \begin{cases} \text{reject } H_0 & \text{if } \tau^* > q_\alpha \\ \text{accept } H_0 & \text{otherwise} \end{cases}$$

SAMPLING DISTRIBUTION ESTIMATION

Statistical significance by permutation test:

- ▶ Calculate τ_{app}^* and record it as E_0 .
- ▶ For $j = 1 : B$, block-wise randomly shuffle $\{X_i\}_{i \in \mathbb{I}}$ with block size D . Calculate τ_{app}^* on shuffled data and recorded it as E_j .
- ▶ P -value: $\#\{E_j > E_0\}/B$



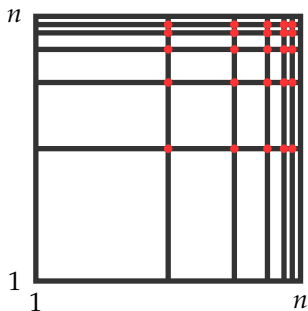
COMPUTATIONAL CONSIDERATION

Fast computation:

$$\tau_f^* := \max_{T_x=X_{(j)}, T_y=Y_{(k)}: j, k \in \mathcal{R}_n} \widehat{\tau}(T)$$

where

$$\mathcal{R}_n := \left\{ s : s = \left\lfloor n - \left(1 + \frac{1}{\log \log n}\right)^j \right\rfloor, j = 1, 2, \dots \quad \text{and} \quad s \geq \lfloor n/2 \rfloor \right\}.$$



- τ_f^* is faster to compute

| | | |
|------------------------|----------|---------------|
| | τ^* | τ_f^* |
| $\# \widehat{\tau}(T)$ | $O(n^2)$ | $O(\log^3 n)$ |

- And just as powerful

DOES IT WORK?

- ▶ Assume that $\{(X_i, Y_i) : i \in \mathbb{I}\}$ ($n := |\mathbb{I}|$) are *independently* sampled from F obeying

$$\sup_{\eta_X, \eta_Y} V(\eta_X, \eta_Y) \cdot T^2(\eta_X, \eta_Y) \gg \frac{\log \log n}{n}.$$

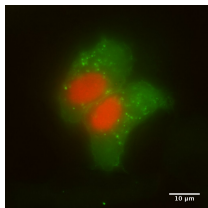
Then Δ is a consistent test in that we reject H_0 in favor of H_1 with probability tending to one.

- ▶ Conversely, there exists a constant $c > 0$ such that for any α -level test Δ based on sample $\{(X_i, Y_i) : i \in \mathbb{I}\}$, there is an instance where joint distribution function F obeying

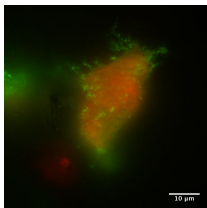
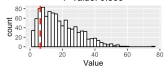
$$\sup_{\eta_X, \eta_Y} V(\eta_X, \eta_Y) \cdot T^2(\eta_X, \eta_Y) \geq c \frac{\log \log n}{n}$$

and yet, we accept H_0 with probability tending to $1 - \alpha$ as if H_0 holds.

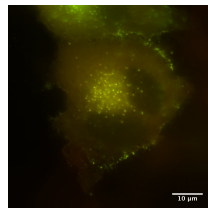
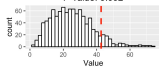
DOES IT REALLY WORK?



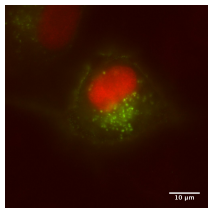
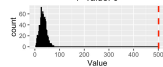
P-value: 0.866



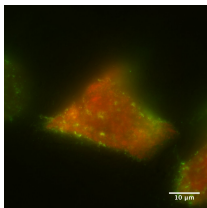
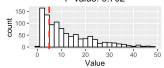
P-value: 0.092



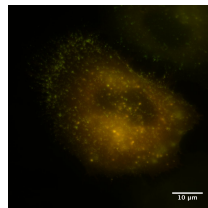
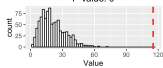
P-value: 0



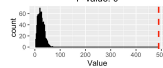
P-value: 0.702



P-value: 0



P-value: 0



What is Colocalization?

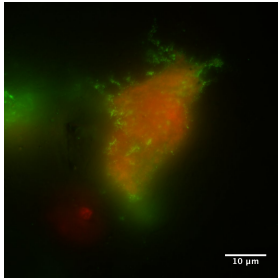
A Statistical View of Colocalization

Global Assessment of Colocalization

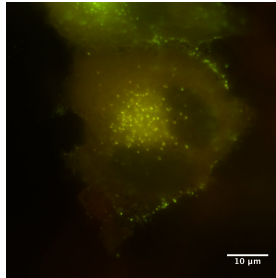
Local Identification of Colocalization

Concluding Remarks

WHERE IS COLOCALIZATION?



$$p = 0.092$$



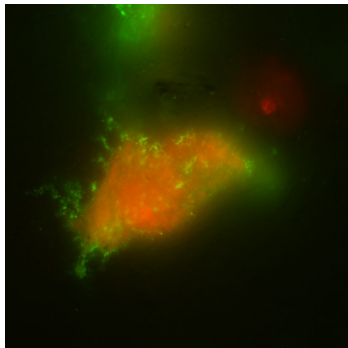
$$p = 0$$

LOCAL QUANTIFICATION OF COLOCALIZATION

- ▶ Pixel-wise hypothesis:

$$H_{k,0} : F_k \in \mathcal{F}_0 \quad \text{v.s.} \quad H_{k,1} : F_k \in \mathcal{F}_1, \quad k \in \mathbb{I}$$

where F_k is the distribution of (X_k, Y_k) .



- ▶ Colocalization as tail dependence:

$$H_{k,0} : Q(F_k; \eta_X, \eta_Y) = 0$$

and

$$H_{k,1} : Q(F_k; \eta_X, \eta_Y) > 0.$$

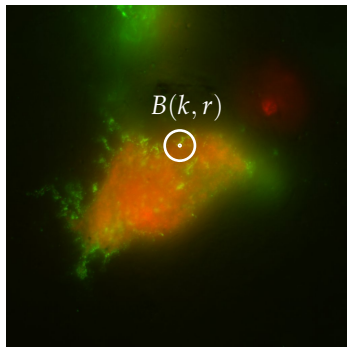
for pre-specified η_X and η_Y .

- ▶ Only one pair (X_k, Y_k) available for each pixel k .

LOCAL QUANTIFICATION OF COLOCALIZATION

Weighted Kendall tau's correlation in a neighborhood $B(k, r)$

$$\tau_w(k; r) := \frac{\sum_{i \neq j} w_i(k; r) w_j(k; r) \text{sign}(X_i - X_j) \text{sign}(Y_i - Y_j)}{\sum_{i \neq j} w_i(k; r) w_j(k; r)}$$



- ▶ Weight $w_i(k; r)$ is decomposed as

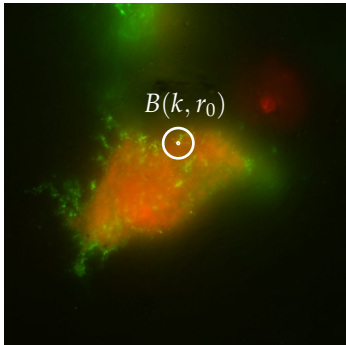
$$w_i(k; r) = \mathbf{K}_l \left(\frac{d(i, k)}{r} \right) \mathbf{K}_b(X_i, Y_i)$$

- ▶ \mathbf{K}_l gives less weight to the pixel i whose location is far from k .
- ▶ $\mathbf{K}_b(X_i, Y_i) = \mathbf{1}_{(X_i > \eta_X)} \mathbf{1}_{(Y_i > \eta_Y)}$ deals with background.

$$w_i(k; r_0)$$

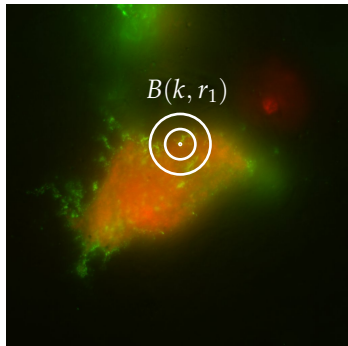


$$\tau_w(k; r_0)$$



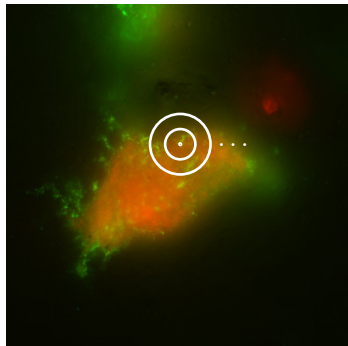
PROPAGATION-SEPARATION

$$\begin{array}{ccc} w_i(k; r_0) & & w_i(k; r_1) \\ \downarrow & \nearrow & \downarrow \\ \tau_w(k; r_0) & & \tau_w(k; r_1) \end{array}$$

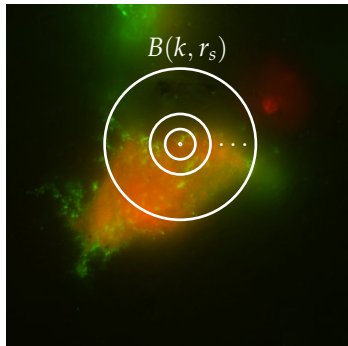
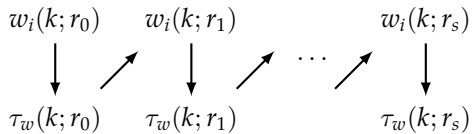


PROPAGATION-SEPARATION

$$\begin{array}{ccc} w_i(k; r_0) & w_i(k; r_1) & \dots \\ \downarrow \nearrow & \downarrow \nearrow & \\ \tau_w(k; r_0) & \tau_w(k; r_1) & \dots \end{array}$$

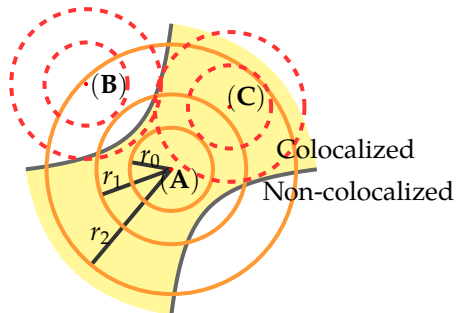
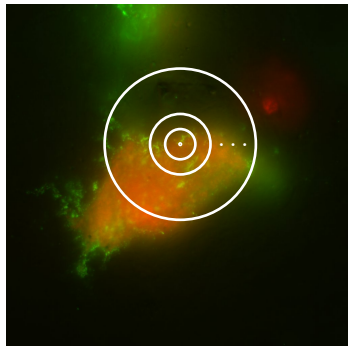


PROPAGATION-SEPARATION



PROPAGATION-SEPARATION

$$\begin{array}{ccccccc} w_i(k; r_0) & & w_i(k; r_1) & & \dots & & w_i(k; r_s) \\ \downarrow & \nearrow & \downarrow & \nearrow & & \nearrow & \downarrow \\ \tau_w(k; r_0) & & \tau_w(k; r_1) & & & & \tau_w(k; r_s) \end{array}$$



- ▶ Test statistics for $H_{k,0}$ against $H_{k,1}$ is

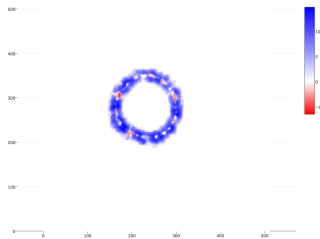
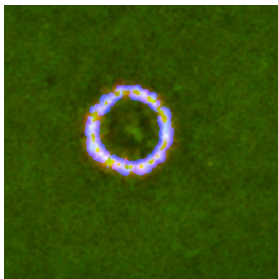
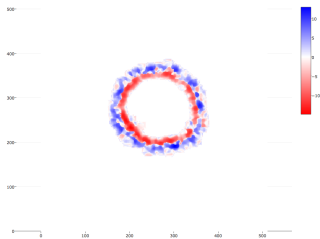
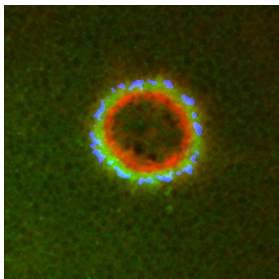
$$Z(k; r_T) = \frac{3}{2} \sqrt{\tilde{N}_k^{(T)}} \cdot \tau_w(k; r_T)$$

where

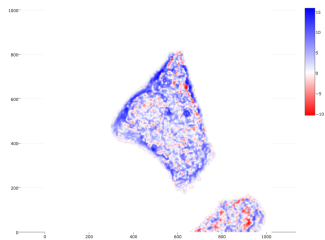
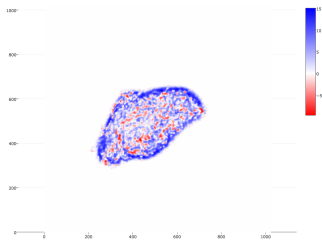
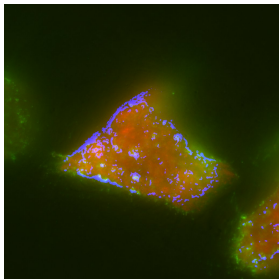
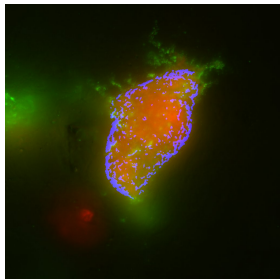
$$\tilde{N}_k^{(T)} = \left(\sum_i w_i(k; r_T) \right)^2 / \sum_i w_i^2(k; r_T).$$

- ▶ Under $H_{k,0}$ s, $Z(k; r_T)$ s behave like standard normal distribution.
- ▶ Correct for multiple testing issue by either Bonferroni method, false discovery rate method or random field theory.

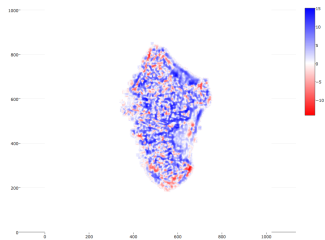
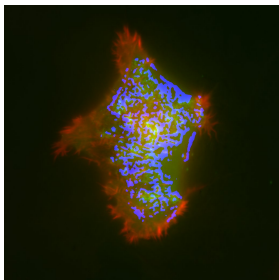
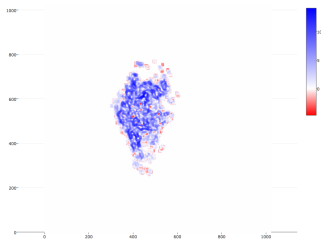
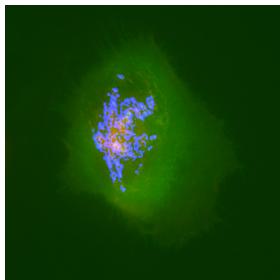
EXAMPLE



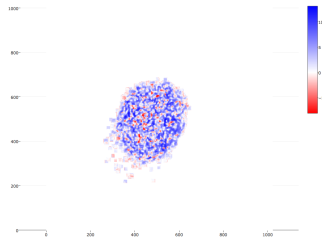
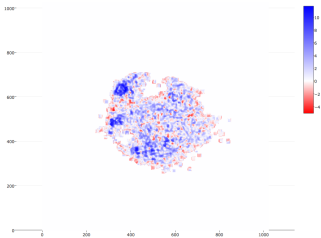
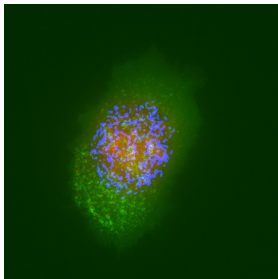
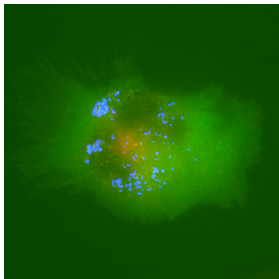
EXAMPLE



EXAMPLE



EXAMPLE



What is Colocalization?

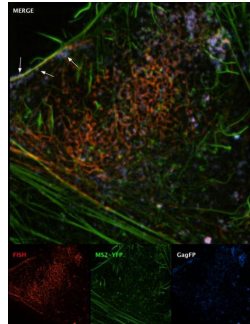
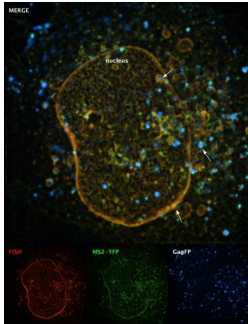
A Statistical View of Colocalization

Global Assessment of Colocalization

Local Identification of Colocalization

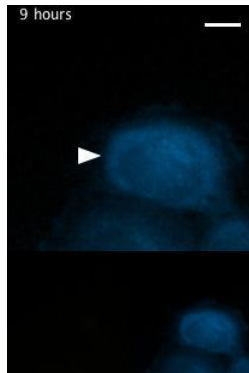
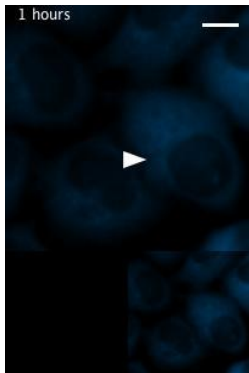
Concluding Remarks

MULTIPLE CHANNELS



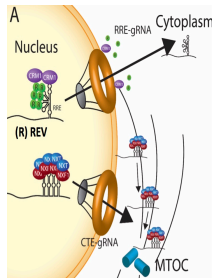
(Becker and Sherer, 2017)

DYNAMICAL COLOCALIZATION



(Becker and Sherer, 2017)

COLLABORATIVE TEAM



- ▶ Colocalization analysis is wide used
- ▶ Quantitation in colocalization analysis
- ▶ Challenges in quantitative imaging

Thank you!