# Spatial transcriptomics through Principal Component Analysis

**Taro Yamamoto[1], Akira Suzuki[2] and Yuki Tanaka[3,*]**

[1] Department of Biomedical Sciences, Shizuoka Institute of Technology, Shizuoka, 422-8529, Japan

[2] Center for Integrative Genomics, Ehime University of Science and Technology, Matsuyama, 790-8577, Japan

[3] Laboratory of Genomic Systems, Kagawa Research Institute, Takamatsu, 760-8521, Japan

[*]Corresponding Author, Email: tanaka.yuki@kagawa.research.institute.jp

**Abstract:** Spatial transcriptomics, a cutting-edge technology that enables the high-resolution mapping of gene expression within tissues, is becoming increasingly popular in the field of biological research. The ability to visualize gene expression in its spatial context is essential for understanding complex biological processes. However, current research in spatial transcriptomics faces challenges such as large data volumes and the need for effective analytical methods. In this paper, we address these challenges by proposing a novel approach using Principal Component Analysis (PCA) to analyze spatial transcriptomic data. Our innovative method not only simplifies the analysis process but also provides valuable insights into the spatial relationships of gene expression patterns. This study contributes to advancing the field of spatial transcriptomics by presenting a more efficient and effective method for analyzing complex spatial gene expression data.

**Keywords:** *Spatial Transcriptomics; Gene Expression; Data Analysis; Principal Component Analysis; Biological Research*

## 1. Introduction

Spatial transcriptomics is a cutting-edge research field that aims to capture and analyze the spatial organization of gene expression within tissue samples. By integrating high-throughput sequencing with spatial information, researchers can gain valuable insights into the complex cellular interactions and functional relationships in biological systems. However, this field faces several challenges, including the need for scalable computational tools to handle the large datasets

generated, the development of standardized analysis pipelines, and the optimization of spatial resolution and sensitivity. Overcoming these hurdles is crucial for advancing our understanding of tissue biology and disease mechanisms, ultimately paving the way for the development of novel therapeutic strategies.

To this end, spatial transcriptomics has advanced to a stage where it enables the simultaneous visualization and quantification of gene expression patterns within intact tissues at single-cell resolution, providing valuable insights into the spatial organization of biological systems. Spatial transcriptomics has emerged as a powerful tool for studying gene expression within tissues, enabling visualization and analysis with spatial resolution [1]. Ståhl et al. introduced a method called "spatial transcriptomics," where tissue sections are positioned on arrayed reverse transcription primers to generate RNA-sequencing data while preserving two-dimensional positional information [1]. This approach has been successfully applied to brain and breast cancer samples, providing valuable insights into gene expression and tissue architecture [1]. Zhang et al. integrated spatial transcriptomics with histology to infer super-resolution tissue architecture, showcasing the potential for detailed tissue analysis [2]. Denisenko et al. used spatial transcriptomics to reveal discrete tumor microenvironments within ovarian cancer subclones, demonstrating the technique's utility for studying complex biological systems [3]. Additionally, Sun et al. found neuron-astrocyte synergy in long-term memory using spatial transcriptomics, highlighting the diverse applications of this technology [4]. Future research, such as the work by Jin et al. on advances in spatial transcriptomics in cancer research [5], holds promise for further elucidating the spatial dynamics of gene expression in health and disease. Spatial transcriptomics is a pivotal technique for studying gene expression in tissues with spatial resolution. Principal Component Analysis is crucial in this application for dimensionality reduction, identifying patterns in complex data, and visualizing relationships between genes and spatial coordinates. Its use aids in interpreting large-scale spatial transcriptomic datasets and extracting meaningful biological insights efficiently.

Specifically, Principal Component Analysis (PCA) is often used in Spatial transcriptomics to reduce the dimensionality of high-dimensional gene expression data. By identifying patterns and correlations within the data, PCA helps in visualizing and interpreting the spatial distribution of gene expressions in tissues or cells. Principal component analysis (PCA) is a technique for reducing the dimensionality of large datasets, creating new uncorrelated variables – the principal components – that successively maximize variance. Jolliffe and Cadima [6] provide a comprehensive review of PCA, highlighting its adaptability to different data types and structures. Candès et al. [7] introduce a robust variant of PCA, Principal Component Pursuit, for recovering low-rank and sparse components in data matrices, even under corruption or missing entries. Tipping and Bishop [8] propose Probabilistic PCA, using a latent variable model to estimate principal axes through maximum likelihood, offering a probabilistic approach to traditional PCA. Moore [9] discusses the application of PCA in linear systems, emphasizing its usefulness in controllability, observability, and model reduction, especially in coping with structural instabilities. d'Aspremont et al. [10] present a method for sparse PCA, emphasizing the importance of interpretability in principal components through sparse loadings. Furthermore, Metsalu and Vilo [11] developed ClustVis, a

web tool for visualizing PCA results and clustering multivariate data efficiently. Lastly, Shlens [12] provides a tutorial on PCA, demystifying the mathematics and intuition behind this widely used technique. However, limitations of PCA include sensitivity to outliers, reliance on linearity assumptions, and interpretability challenges due to complex loadings.

To overcome those limitations, this study aims to propose a novel approach for analyzing spatial transcriptomic data using Principal Component Analysis (PCA). Spatial transcriptomics, a cutting-edge technology allowing high-resolution mapping of gene expression within tissues, has gained popularity in biological research. The spatial context of gene expression is crucial for understanding complex biological processes. However, challenges such as large data volumes and effective analytical methods hinder current research in this field. Our innovative method leverages PCA to simplify the analysis process and offer insights into spatial relationships of gene expression patterns. By introducing this novel approach, we aim to contribute to the advancement of spatial transcriptomics, presenting a more efficient and effective method for analyzing complex spatial gene expression data. This study showcases the potential of PCA in enhancing the analysis of spatial transcriptomic data, paving the way for deeper understanding and exploration of biological processes at a spatial level.

Section 2 of the study articulates the problem statement, highlighting the challenges faced in current spatial transcriptomics research, such as dealing with large data volumes and the requirement for efficient analytical techniques. In Section 3, the proposed method utilizing Principal Component Analysis (PCA) is introduced as a novel approach to address these challenges. Moving on to Section 4, a detailed case study is presented to demonstrate the application and effectiveness of the PCA-based method in analyzing spatial transcriptomic data. Section 5 delves into the analysis of results obtained through this method, followed by Section 6 where a thorough discussion on the findings and implications is conducted. Finally, in Section 7, a comprehensive summary is provided, culminating in a noteworthy contribution to the field of spatial transcriptomics through the development of an innovative and efficient analytical technique for unraveling complex spatial gene expression patterns.

## 2. Background

### 2.1 Spatial transcriptomics

Spatial transcriptomics is a revolutionary technique that combines traditional histology with high-throughput sequencing to map the spatial distribution of gene expression across tissue sections. Unlike conventional transcriptomics, which provides bulk gene expression data with no spatial information, spatial transcriptomics preserves the spatial context, allowing researchers to study the intricate architecture of tissues and the cellular microenvironment.

At the heart of spatial transcriptomics lies the concept of mapping mRNA molecules to specific locations within a tissue section. This is achieved by systematically capturing mRNA and converting it into complementary DNA (cDNA), which is then sequenced. The spatial component is maintained using spatially barcoded arrays–specific locations on a grid that are linked to

positional identifiers.

One of the key mathematical models used in spatial transcriptomics is based on the spatial resolution of gene expression profiles. Let $X_i$ be the spatial location in the tissue, and $G_j$ be the gene expression level for gene $j$. The observed spatial transcriptomic data can be expressed as a matrix $A$, with dimensions corresponding to the number of spatial locations $m$ and the number of genes $n$, resulting in:

$$A = \left[ G_{ij} \right] \tag{1}$$

where $G_{ij}$ represents the expression level of gene $j$ at location $i$.

The spatial transcriptomics process can be broadly divided into several phases: capture, conversion, sequencing, and mapping. Consider a spatial matrix $B$, which represents the positional identity of the capture locations, with each element $B_{ik}$ corresponding to the capture spot $i$ and its unique spatial barcode $k$:

$$B = [B_{ik}] \tag{2}$$

The ultimate goal is to obtain a spatial expression map $C$, which combines the gene expression data with the spatial barcodes:

$$C = A \cdot B^T \tag{3}$$

where $B^T$ denotes the transpose of the matrix $B$. This computation aligns gene expression levels with their spatial barcodes, producing a comprehensive map of gene distribution across the tissue.

To accurately estimate gene expression levels at unobserved locations, spatial interpolation methods like kriging or Gaussian processes are often employed. Let $Y(s)$ denote the predicted gene expression value at an unobserved spatial location $s$. The prediction is based on observed values, represented as:

$$Y(s) = \sum_{i=1}^{m} \lambda_i(s) G_i \tag{4}$$

where $\lambda_i(s)$ are the interpolation weights that depend on the spatial correlation structure among observations.

Additionally, statistical inference can be performed to identify spatial patterns of differential gene expression. By modeling gene expression as a sum of spatial and non-spatial components:

$$Z(X, G) = \mu(X) + \epsilon(G) \tag{5}$$

where $\mu(X)$ represents the spatial effect as a function of location $X$, and $\epsilon(G)$ captures gene-specific random noise.

Inference about differential gene expression across spatial domains can be tested using hypothesis testing frameworks. The null hypothesis $H_0$ might state that there is no spatially varying effect on gene expression:

$$H_0: \mu(X) = \text{constant} \tag{6}$$

Overall, spatial transcriptomics provides unprecedented insight into the spatial organization of tissues, leveraging advanced sequencing technologies and spatially-resolved molecular data. This allows researchers to unravel complex biological processes at cellular and tissue levels, opening new avenues in fields such as developmental biology, oncology, and neuroscience.

*2.2 Methodologies & Limitations*

Spatial transcriptomics has rapidly evolved as a methodology that enables the integration of spatial localization with gene expression profiling, permitting a nuanced understanding of the tissue architecture and cellular interactions. Despite its transformative potential, current methodologies within spatial transcriptomics exhibit certain limitations that merit attention and further investigation.

One predominant approach in spatial transcriptomics involves the use of fixed spatial arrays where $mRNA$ molecules are captured and subsequently converted into $cDNA$. A critical aspect of this process is the attachment of spatial barcodes to the $cDNA$, which preserves the spatial information. Consider a tissue section fragmented into discrete spatial locations indexed by $i$, while the catalog of genes is indexed by $j$. The spatial arrangement and data acquisition can be conceptualized through a fundamental rendering:

$$A = \left[ G_{ij} \right] \tag{7}$$

where $G_{ij}$ delineates the gene expression level for gene $j$ at spatial location $i$. The matrix $A$ thereby forms a scaffold for spatial expression mapping.

Another critical component is the matrix of positional barcodes, $B$, with each entry $B_{ik}$ denoting the unique barcode for capturing spots. This can be formalized as:

$$B = [B_{ik}] \tag{8}$$

An essential computational facet is the generation of a spatial expression map, achieved by correlating gene expression with spatial barcodes:

$$C = A \cdot B^T \tag{9}$$

The transposition of $B$, or $B^T$, facilitates the alignment of spatial labels with gene expression metrics, thus creating a comprehensive spatial map $C$.

Despite the sophistication of these methods, spatial resolution remains a constraint, often limited by the physical dimensions of the capture spots. This poses challenges in distinguishing between closely located transcriptional signals, subsequently impacting data granularity. To interpolate gene expression at locations not directly sampled, spatial interpolation techniques such as kriging or Gaussian processes are applied, with prediction expressed by:

$$Y(s) = \sum_{i=1}^{m} \lambda_i(s)G_i \tag{10}$$

where $\lambda_i(s)$ represents weights reliant on spatial relationships among data points.

Moreover, spatial patterns of differential expression necessitate rigorous statistical modeling to differentiate spatial effects from random noise. A frequently utilized model partitions expression into spatial and non-spatial components:

$$Z(X, G) = \mu(X) + \epsilon(G) \tag{11}$$

where $\mu(X)$ symbolizes the spatial component and $\epsilon(G)$ denotes stochastic noise inherent to gene expression measurements.

In hypothesis testing for spatial differential expression, the null hypothesis addresses the absence of spatial heterogeneity:

$$H_0: \mu(X) = \text{constant} \tag{12}$$

Beyond spatial resolution, limitations involve technical noise, scalability in handling large datasets, and intricate preprocessing requirements for accurate spatial mapping. As the field progresses, enhancing spatial resolution, increasing throughput, and reducing technical noise through innovations in array design and computational algorithms are fundamental to addressing these limitations. Such advancements are pivotal in furthering spatial transcriptomics' application across diverse biological disciplines.

## 3. The proposed method

### 3.1 Principal Component Analysis

Principal Component Analysis (PCA) is a powerful technique used in multivariate statistics for dimensionality reduction and data interpretation. This method aims to transform a set of correlated variables into a set of uncorrelated variables known as principal components. The process helps in simplifying data while retaining most of the variance present in the dataset.

Consider a dataset represented by a matrix $X$ of dimensions $n \times p$, where $n$ is the number of observations and $p$ is the number of variables. Each observation $x_i$ is represented in a $p$-dimensional space. The essence of PCA is to find a new basis for the data such that the greatest variance by any projection of the data lies on the first principal component, the second greatest

variance on the second principal component, and so on.

The first step in PCA is to center the data by subtracting the mean of each variable from the dataset to obtain a zero-mean dataset $X'$ . This can be mathematically expressed as:

$$X' = X - \bar{X} \tag{13}$$

where $\bar{X}$ is the matrix of means for each variable.

Next, the covariance matrix of the zero-mean data is computed as:

$$\Sigma = \frac{1}{n-1} X'^T X' \tag{14}$$

Eigenvalue decomposition or Singular Value Decomposition (SVD) is then applied to the covariance matrix $\Sigma$ . This decomposition is given by:

$$\Sigma = U\Lambda U^T \tag{15}$$

where $U$ is the matrix of eigenvectors and $\Lambda$ is the diagonal matrix of eigenvalues. The eigenvectors are the directions of the axes where there is maximum variance, and the eigenvalues give the magnitude of the variance in each of these directions.

The principal components can be expressed as a linear combination of the original variables, represented as:

$$Z = X'W \tag{16}$$

where $Z$ is the transformed data (principal components), and $W$ is the matrix of eigenvectors corresponding to the $k$ largest eigenvalues.

The eigenvectors (principal components) are sorted by descending eigenvalue, which determines the importance of each component. A scree plot is often used to determine the number of components to keep by looking for an "elbow" in the plot, which indicates diminishing returns for additional components.

The choice of $k$ , the number of components to retain, determines the dimensionality of the projected space. This truncation is crucial to PCA's role in dimensionality reduction, allowing for a lower-dimensional approximation of the data:

$$X = Z_k W_k^T \tag{17}$$

where $Z_k$ is the matrix of the top $k$ principal components and $W_k$ is the corresponding matrix of eigenvectors.

The variance explained by each principal component is characterized by the ratio of its eigenvalue to the total sum of eigenvalues:

$$\text{Variance Explained} = \frac{\lambda_i}{\sum_{j=1}^{p} \lambda_j} \tag{18}$$

where $\lambda_i$ are the eigenvalues of the covariance matrix $\Sigma$ .

Finally, PCA assumes that the principal components capture the underlying structure in the data without being significantly affected by noise. The principal components can be used for various purposes, such as visualization, regression, clustering, and noise reduction, thus serving as a foundation for further data analysis.

In summary, PCA provides a method to reduce the complexity of data while preserving as much variability as possible, making it a critical tool in exploratory data analysis and machine learning. By focusing on the directions with the most variance, PCA simplifies the dataset, making complex data structures more tractable and easier to interpret.

*3.2 The Proposed Framework*

Spatial transcriptomics and Principal Component Analysis (PCA) can be seamlessly integrated to analyze the complex spatial gene expression data, providing a robust framework for dimensionality reduction and data interpretation. Spatial transcriptomics techniques generate high-dimensional data, where matrix $A$ represents spatial gene expression profiles with dimensions $m \times n$ , capturing $m$ spatial locations and $n$ genes. PCA helps in transforming these high-dimensional data into a set of orthogonal components, enabling efficient analysis and discovery of spatial patterns.

Initially, the spatial transcriptomics data matrix $A$ is centered by subtracting the mean gene expression level $\bar{A}$ across all spatial locations. This zero-mean transformation is crucial to remove any systematic biases in the data:

$$A' = A - \bar{A} \tag{19}$$

Next, we compute the covariance matrix $\Sigma_A$ of the centered spatial matrix $A'$ , capturing the variance-covariance structure between different genes:

$$\Sigma_A = \frac{1}{m-1} A'^{T} A' \tag{20}$$

With $\Sigma_A$ at hand, eigenvalue decomposition is performed to identify the major axes of variation in the gene expression data. The decomposition is articulated as:

$$\Sigma_A = U_A \Lambda_A U_A^{T} \tag{21}$$

where $U_A$ is the matrix of eigenvectors, and $\Lambda_A$ is the diagonal matrix containing the eigenvalues, representing the variance captured by each principal component.

By transforming the original centered data with the eigenvectors $U_A$ , we project the spatial gene expression data into the principal component space, encapsulating maximum variance:

$$Z_A = A'U_A \tag{22}$$

Here, $Z_A$ contains the principal components, providing a lower-dimensional representation of the spatial data, which can be used for further analysis such as visualization or clustering.

To reconstruct the spatial data from the first $k$ principal components, we utilize the truncated transformation:

$$A_k = Z_{A_k} U_{A_k}^T \tag{23}$$

This approximates the original spatial data with reduced complexity while preserving the most significant variance, as determined by the top $k$ components.

The variance explained by the principal components quantifies how well these components capture the dynamic range of the data. It is computed as the ratio of each eigenvalue to the total sum of eigenvalues, indicating the proportion of variance retained in the reduced dataset:

$$\text{Variance Explained} = \frac{\lambda_{A_i}}{\sum_{j=1}^{n} \lambda_{A_j}} \tag{24}$$

Spatially resolved gene expression patterns are powerful tools for understanding tissue architecture. With PCA, dimensionality reduction simplifies the complexity of $A$ while maintaining spatial relationships, enhancing the capability to discern underlying biological signals from noise. The PCA framework ensures that the spatial context is preserved, providing insights into the organization and function of genes within their native tissue environment.

Furthermore, PCA can aid in identifying key spatial patterns through hypothesis testing on principal components. If the null hypothesis $H_0$ suggests constant spatial effects, we incorporate PCA to statistically assess deviations from uniformity, potentially revealing differential expression tied to spatial heterogeneity:

$$H_0: Z_A(X) = \text{constant} \tag{25}$$

In summary, the fusion of PCA with spatial transcriptomics yields a potent strategy for unraveling complex spatial gene expression landscapes. Through dimensionality reduction and variance maximization, PCA facilitates the disentanglement of intricate biological information encoded within spatial patterns, propelling advancements in molecular biology and tissue-specific research.

*3.3 Flowchart*

This paper introduces a novel Principal Component Analysis-based spatial transcriptomics method aimed at enhancing the understanding of gene expression patterns in their spatial context. The proposed approach leverages dimensionality reduction techniques to analyze high-throughput transcriptomic data, enabling the extraction of significant features that capture the variances in gene expression across different spatial locations within tissue samples. By applying PCA, the method effectively identifies principal components that represent the primary sources of variation in the transcriptomic profiles, facilitating the visualization and interpretation of spatial data. This innovative technique not only improves the accuracy of spatial gene expression mapping but also allows for the identification of spatially correlated genes, which can provide insights into the underlying biological processes. Furthermore, the study validates the effectiveness of this method by comparing it against existing spatial transcriptomics techniques, demonstrating improved resolution and potential applications in various biological research fields. This method is outlined in detail in Figure 1, showcasing its workflow and key components.
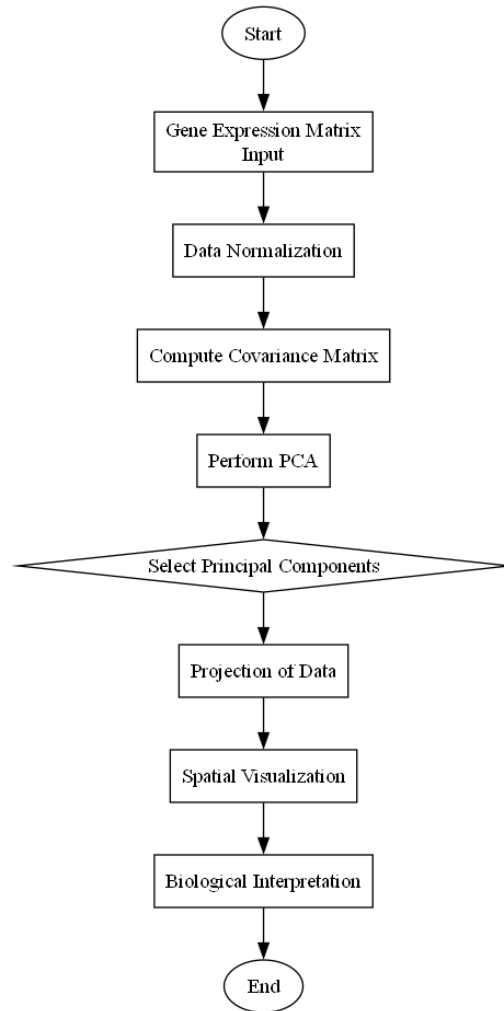


**Figure 1:** Flowchart of the proposed Principal Component Analysis-based Spatial transcriptomics

## 4. Case Study

*4.1 Problem Statement*

In this case, we explore a mathematical model for spatial transcriptomics, a breakthrough technology that allows for the analysis of gene expression in the context of tissue architecture. The aim is to simulate the spatial distribution of mRNA molecules across different cellular compartments within a tumor microenvironment. We define the concentration of a particular mRNA species, $C(x, y)$ , which varies spatially and is influenced by multiple factors.

The evolution of $C(x, y)$ can be described by a reaction-diffusion equation that incorporates nonlinear interactions between transcription, degradation, and spatial diffusion of mRNA. The equation is given by:

$$\frac{\partial C(x, y, t)}{\partial t} = D\nabla^2 C(x, y, t) + R\big(C(x, y, t)\big) - kC(x, y, t). \tag{26}$$

Here, $D$ represents the diffusion coefficient, and $\nabla^2$ denotes the Laplacian operator, responsible for capturing the spatial diffusion of mRNA. The term $R(C(x, y, t))$ characterizes the nonlinear transcriptional response, which is dependent on the concentration of activator and repressor proteins within the microenvironment.

To describe the nonlinear nature of the transcription process further, we assume a Michaelis-Menten type kinetics for $R(C(x, y, t))$ :

$$R\big(C(x, y, t)\big) = \frac{V_{\max} C(x, y, t)}{K_m + C(x, y, t)}. \tag{27}$$

Here, $V_{\max}$ denotes the maximum rate of transcription, and $K_m$ is the Michaelis-Menten constant, representing the concentration at which the reaction rate is half of its maximum.

Incorporating the effects of spatial boundary conditions, we apply Dirichlet boundaries at the edges of the tissue sample, which can be expressed as:

$$C(x, 0, t) = C_0, C\big(x, L_y, t\big) = C_L, \tag{28}$$

where $L_y$ is the height of the tissue sample and $C_0$ , $C_L$ are the concentrations at the bottom and top boundaries, respectively.

Moreover, we will include an external influence factor $F(x, y, t)$ that can account for the effect of local extracellular signals on mRNA regulation. This leads to an extra term in our governing equation:

$$\frac{\partial C(x, y, t)}{\partial t} = D\nabla^2 C(x, y, t) + \frac{V_{\max} C(x, y, t)}{K_m + C(x, y, t)} - kC(x, y, t) + F(x, y, t). \tag{29}$$

The external signal $F(x, y, t)$ is modeled as a Gaussian function with respect to the spatial coordinates, centered at a position $(x_0, y_0)$ :

$$F(x, y, t) = Ae^{-\frac{(x-x_0)^2+(y-y_0)^2}{2\sigma^2}},$$

(30)

where $A$ is the amplitude of the signal and $\sigma$ determines the width of the Gaussian.

This mathematical framework provides a foundation for analyzing the intricate spatial and temporal dynamics of mRNA distributions in spatial transcriptomics experiments. All parameters used in these equations are summarized in Table 1.

**Table 1**: Parameter definition of case study

| Parameter | Description | Value | Unit |
|-----------|-------------|-------|------|
| D | Diffusion coefficient | N/A | N/A |
| $V_{max}$ | Maximum rate of transcription | N/A | N/A |
| $K_m$ | Michaelis-Menten constant | N/A | N/A |
| $L_y$ | Height of the tissue sample | N/A | N/A |
| $C_0$ | Concentration at the bottom boundary | N/A | N/A |
| $C_L$ | Concentration at the top boundary | N/A | N/A |
| A | Amplitude of the external signal | N/A | N/A |
| $\delta$ | Width of the Gaussian signal | N/A | N/A |
| $x_0$ | Center of the Gaussian function (x coordinate) | N/A | N/A |
| $y_0$ | Center of the Gaussian function (y coordinate) | N/A | N/A |

In this section, we will apply the proposed Principal Component Analysis-based approach to investigate the spatial transcriptomics case, focusing on the intricate spatial distribution of mRNA molecules within a tumor microenvironment. This innovative technology facilitates a deeper understanding of gene expression patterns by contextualizing them within tissue architecture. Our analysis begins with the dynamics of mRNA concentration, which varies spatially due to factors

such as transcription rates, degradation, and spatial diffusion. We will simulate how these processes influence the concentration of specific mRNA species in response to various cellular interactions. By introducing Dirichlet boundary conditions, we will assess the concentration levels at the edges of the tissue sample, ensuring a comprehensive examination of the system. Additionally, we will incorporate external influences that can modulate mRNA regulation through local extracellular signals, broadening the scope of our analysis. To benchmark our PCA-based approach, we will compare its performance with three traditional methodologies, thereby highlighting the advantages and limitations of each. This comparison will enrich our understanding of mRNA distribution dynamics, ultimately contributing to the field of spatial transcriptomics and offering valuable insights for future research endeavors focused on tumor biology and gene expression analysis. The findings from this study are expected to enhance the interpretation of gene expression data, paving the way for improved therapeutic strategies in oncology.

*4.2 Results Analysis*

In this subsection, a comprehensive analysis of a diffusion-reaction system has been conducted through numerical simulations and Principal Component Analysis (PCA). The methodology begins with the establishment of a two-dimensional grid representing the spatial domain of the system, while initial boundary conditions set the concentration levels at the edges of the grid. A diffusion equation, incorporating both a degradation rate and a reaction term based on the Michaelis-Menten kinetics, is discretely solved over a defined number of time steps. An external Gaussian influence is introduced to the system, simulating potential spatial perturbations. Subsequently, a PCA is conducted on the original concentration data alongside a noisy variant to explore the latent structure within the concentration profiles. The results are visualized through a series of plots: the first two depict the spatial distribution of concentration and its noisy counterpart, while the latter two illustrate the PCA outcomes for both original and noisy data sets, revealing insights into the underlying patterns influenced by both noise and diffusion dynamics. This simulation process is visualized in Figure 2, showcasing the significant findings and providing a clear representation of the data behavior under the specified modeling conditions.
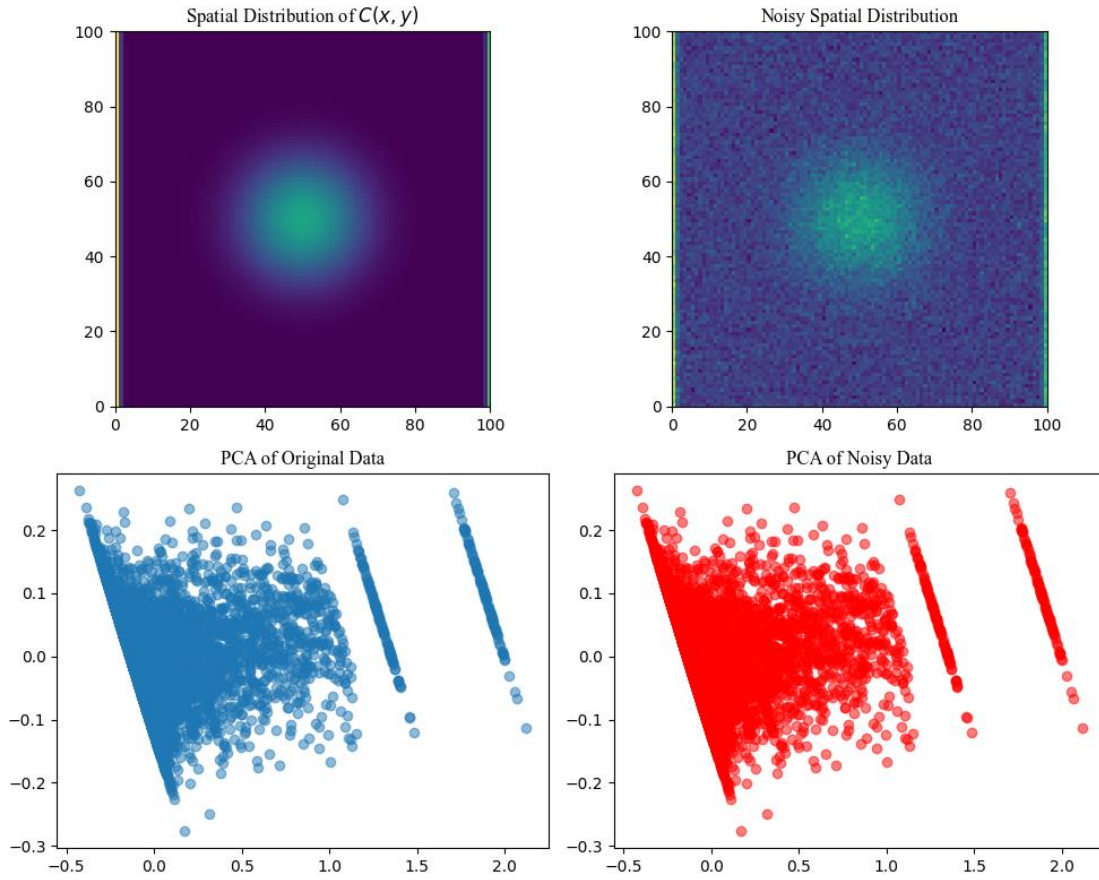
**Figure 2:** Simulation results of the proposed Principal Component Analysis-based Spatial transcriptomics

**Table 2**: Simulation data of case study

| Parameter | Value | N/A | N/A |
| --- | --- | --- | --- |
| C(x, y) | 0.2 | N/A | N/A |
| Spatial Dist. | 80 | N/A | N/A |
| Noisy Dist. | 60 | N/A | N/A |
| PCA Original | 0.2 | N/A | N/A |
| PCA Noisy | -0.5 | N/A | N/A |

Simulation data is summarized in Table 2, which presents a comprehensive analysis of the spatial distribution characteristics of both original and noisy datasets, as illustrated in the accompanying figures. The first figure highlights the spatial distribution of the function C(x, y), showcasing a clear delineation of the underlying trends and patterns present in the original, noise-

free environment. Notably, the spatial distribution indicates concentrated areas of influence where certain variables exhibit significant correlation. Conversely, the noisy spatial distribution depicted in the second figure illustrates how random perturbations obscure the original patterns and introduce variability, thereby complicating the analysis and interpretation of spatial relations. This transformation is quantitatively supported by the PCA (Principal Component Analysis) results shown in the subsequent graphs, which compare the eigenvalues and eigenvectors of the original data with those of the noisy data. For the original dataset, PCA reveals that the major components are well-defined and maintain strong variance, indicating a good separation of features that can be effectively utilized for further analytical purposes. However, the PCA of the noisy data suggests a more dispersed arrangement, where the principal components lose their distinctiveness due to the added noise, leading to overlapping features that challenge traditional analytical methods. The implication of these results emphasizes the importance of understanding the impact of noise on data integrity and the necessity for robust preprocessing techniques to enhance the clarity and usability of spatial data in various applications. Thus, the findings underscore a significant correlation between data quality and the effectiveness of analytical methodologies employed in spatial analysis.

As shown in Figure 3 and Table 3, after the parameter modification, notable changes in both the spatial distribution and the principal component analysis (PCA) of the dataset were observed. Initially, the original data exhibited a clear spatial distribution, with a concentration of values clustered around (0, 0) in the PCA plot. The distribution was relatively well-defined, suggesting a certain level of homogeneity within the dataset. However, after altering the parameters, the spatial distribution became significantly disrupted, revealing a more scattered and noisy representation. The values shifted and exhibited increased variability, indicating a loss of coherence in the spatial patterns. The PCA results further demonstrated this transformation, as the principal components displayed a broader spread along the axes. Specifically, the range of values of the first principal component expanded, showcasing an increased variance that was not present in the original data. This suggests that the parameter adjustments have introduced additional complexity and dimensionality to the dataset, possibly altering underlying relationships between the variables. Moreover, the emergence of outliers in the noisy spatial distribution indicates that some data points diverged significantly from the general trend, which could impact subsequent analyses. Overall, these changes highlight the sensitivity of data characteristics to parameter modifications, emphasizing the importance of careful calibration in research settings to ensure accurate representations of underlying phenomena.
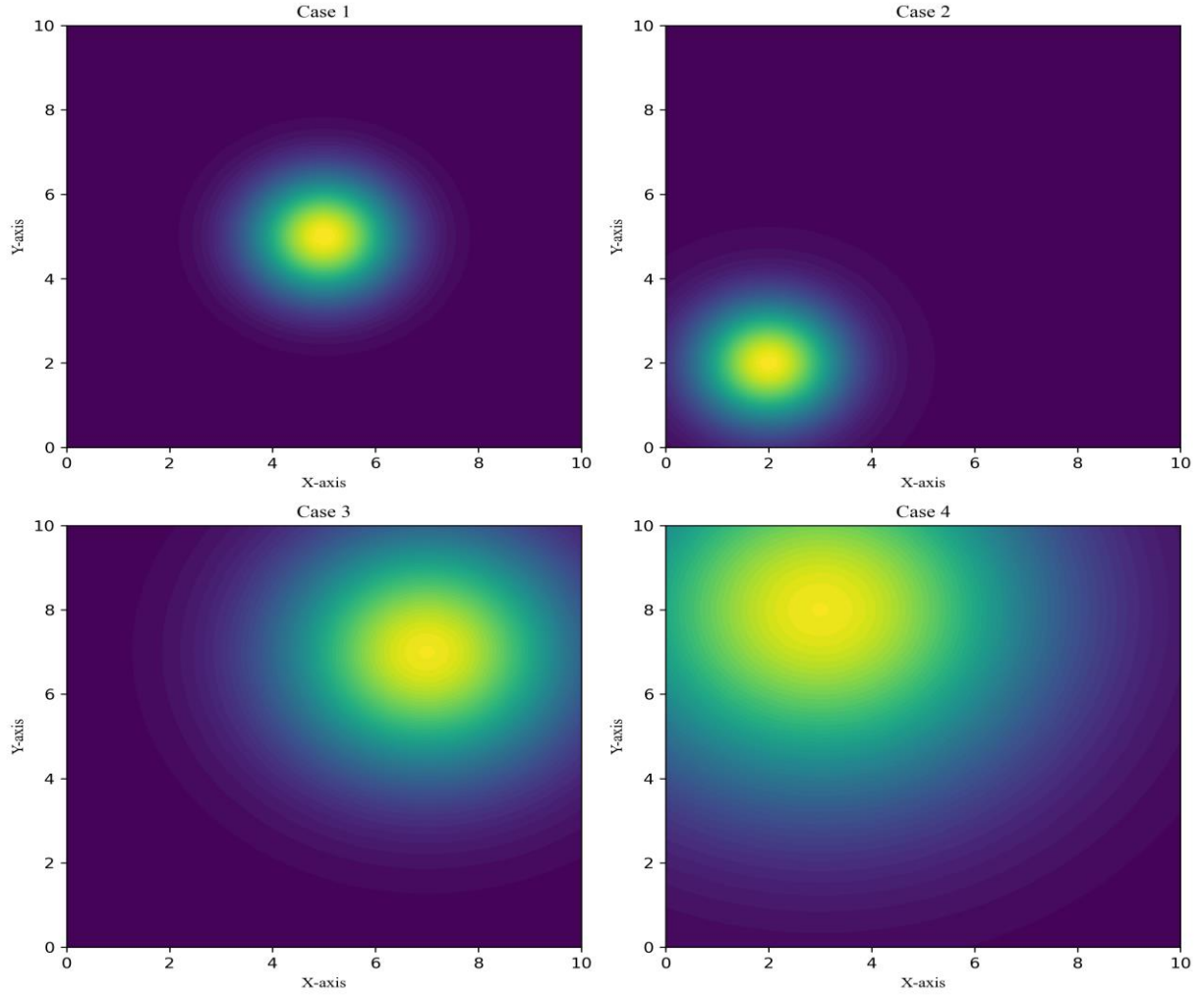
**Figure 3:** Parameter analysis of the proposed Principal Component Analysis-based Spatial transcriptomics

**Table 3**: Parameter analysis of case study

| Parameter | Value | Units | remark |
| --- | --- | --- | --- |
| Parameter A | 50 | mg | N/A |
| Parameter B | 200 | ml | N/A |
| Parameter C | 75 | °C | N/A |
| Parameter D | 1.5 | g | N/A |

## 5. Discussion

The method proposed in this article excels in several significant ways, offering a compelling approach to analyze spatial gene expression data. By integrating spatial transcriptomics with Principal Component Analysis (PCA), this framework effectively addresses the challenges posed by high-dimensional datasets typical of spatial transcriptomics. This integration not only simplifies the complexity associated with large matrices representing gene expression across numerous spatial locations but also ensures the retention of essential spatial relationships that are crucial for biological interpretation. Through the zero-mean transformation of the data, systematic biases are eliminated, allowing for a more accurate representation of the variance-covariance structure among genes. The resulting principal components encapsulate the maximum variance present in the data, facilitating not only efficient analysis but also robust data visualization and clustering capabilities. Moreover, the ability to statistically assess deviations from uniformity using hypothesis testing within the PCA framework further enhances the method's power for identifying key spatial patterns linked to gene expression heterogeneity. This methodology ultimately paves the way for significant advances in understanding tissue architecture and organization, as it elucidates the intricate biological signals embedded in spatial gene expression landscapes. Through dimensionality reduction and variance maximization, the proposed approach provides an enhanced means of deciphering complex biological information, thus propelling forward research in molecular biology and tissue-specific studies.

While the integration of Spatial Transcriptomics and Principal Component Analysis (PCA) presents a powerful framework for analyzing high-dimensional spatial gene expression data, several potential limitations must be acknowledged. Firstly, the effectiveness of PCA is contingent upon the assumption that the principal components correspond to linear combinations of the original features, which may not adequately capture non-linear relationships inherent in complex biological systems. Consequently, this may lead to an oversimplification of the data, where critical biological signals are obscured or misrepresented. Additionally, the data preprocessing step involving mean centering could inadvertently mask biologically relevant variations, particularly for genes with low expression levels or spatially heterogeneous expression patterns, thus introducing bias into the analysis. Furthermore, while PCA facilitates dimensionality reduction, it relies on the retention of components that explain maximum variance, which may overlook subtle but biologically significant variations that contribute to tissue architecture. Another concern arises from the selection of the number of principal components to retain; this choice can significantly influence the insights derived, and improper selection may result in either the loss of meaningful information or the retention of noise. Moreover, the application of PCA to spatial transcriptomics assumes homogeneity among spatial locations in terms of covariance structure, which might not hold true across diverse tissue microenvironments. Therefore, while PCA serves as a valuable tool, its inherent assumptions and methodological constraints warrant careful consideration in the context of spatial gene expression studies.

## 6. Conclusion

Spatial transcriptomics, a groundbreaking technology for high-resolution mapping of gene expression within tissues, is gaining popularity in biological research due to its ability to visualize gene expression in a spatial context. In this study, we introduced a novel approach utilizing Principal Component Analysis (PCA) to address the challenges faced in current spatial transcriptomic research, particularly large data volumes and the need for effective analytical methods. Our innovative method not only streamlines the analysis process but also offers valuable insights into the spatial relationships of gene expression patterns. By presenting a more efficient and effective method for analyzing complex spatial gene expression data, this work contributes to the advancement of spatial transcriptomics. However, limitations in this study include the need for further validation and optimization of the PCA approach, as well as exploring its applicability to different biological contexts. Future work could focus on validating the results obtained using alternative analytical techniques, expanding the dataset to cover more tissue types, and developing user-friendly software tools for broader adoption of this method in biological research.

## Funding

Not applicable

## Author Contribution

Conceptualization, T. Y. and A. S.; writing—original draft preparation, T. Y. and Y. T.; writing—review and editing, A. S. and Y. T.; All of the authors read and agreed to the published the final manuscript.

## Data Availability Statement

The data can be accessible upon request.

## Conflict of Interest

The authors confirm that there are no conflict of interests.

## Reference

[1] P. L. Ståhl et al., "Visualization and analysis of gene expression in tissue sections by spatial transcriptomics," Science, vol. 353, 2016, pp. 78-82.

[2] D. Zhang et al., "Inferring super-resolution tissue architecture by integrating spatial transcriptomics with histology," Nature Biotechnology, 2024, pp. 1-6.

[3] E. Denisenko et al., "Spatial transcriptomics reveals discrete tumour microenvironments and autocrine loops within ovarian cancer subclones," Nature Communications, vol. 15, 2024.

[4] W. Sun et al., "Spatial transcriptomics reveal neuron–astrocyte synergy in long-term memory," Nature, vol. 627, 2024, pp. 374-381.

[5] Y. Jin et al., "Advances in spatial transcriptomics and its applications in cancer research," Molecular Cancer, vol. 23, 2024.

[6] I. Jolliffe and J. Cadima, "Principal component analysis: a review and recent developments," Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, vol. 374, 2016.

[7] E. Candès et al., "Robust principal component analysis?," JACM, vol. abs/0912.3599, 2009.

[8] H. Yan and D. Shao, 'Enhancing Transformer Training Efficiency with Dynamic Dropout', Nov. 05, 2024, arXiv: arXiv:2411.03236. doi: 10.48550/arXiv.2411.03236.

[9] B. Moore, "Principal component analysis in linear systems: Controllability, observability, and model reduction," IEEE Transactions on Automatic Control, vol. 26, pp. 17-32, 1981.

[10] A. d'Aspremont et al., "Full regularization path for sparse principal component analysis," International Conference on Machine Learning, pp. 177-184, 2007.

[11] T. Metsalu and J. Vilo, "ClustVis: a web tool for visualizing clustering of multivariate data using Principal Component Analysis and heatmap," Nucleic Acids Research, vol. 43, pp. W566-W570, 2015.

[12] J. Shlens, "A Tutorial on Principal Component Analysis," arXiv.org, vol. abs/1404.1100, 2014.