
QuST-LLM: INTEGRATING LARGE LANGUAGE MODELS FOR COMPREHENSIVE SPATIAL TRANSCRIPTOMICS ANALYSIS

A PREPRINT

 **Chao Hui Huang**
Pfizer Inc.
La Jolla, CA 92101

July 1, 2024

ABSTRACT

In this paper, we introduce QuST-LLM, an innovative extension of QuPath that utilizes the capabilities of large language models (LLMs) to analyze and interpret spatial transcriptomics (ST) data. In addition to simplifying the intricate and high-dimensional nature of ST data by offering a comprehensive workflow that includes data loading, region selection, gene expression analysis, and functional annotation, QuST-LLM employs LLMs to transform complex ST data into understandable and detailed biological narratives based on gene ontology annotations, thereby significantly improving the interpretability of ST data. Consequently, users can interact with their own ST data using natural language. Hence, QuST-LLM provides researchers with a potent functionality to unravel the spatial and functional complexities of tissues, fostering novel insights and advancements in biomedical research. QuST-LLM is a part of QuST project. The source code is hosted on GitHub and documentation is available at <https://github.com/huangch/qust>.

Keywords Large language model · spatial transcriptomics · gene ontology · knowledge graph · QuPath extension

1 Introduction

Spatial transcriptomics (ST) (Ståhl *et al.* [2016]) has emerged as a transformative technology in the field of genomics, enabling the high-resolution mapping of gene expression across tissue sections. This spatial context is crucial for understanding the cellular architecture and functional organization of tissues. Traditional transcriptomics, which averages gene expression across entire tissues or cell populations, often obscures the spatial heterogeneity and cell-to-cell variability that are fundamental to tissue function and disease progression (Janesick *et al.* [2023], Nature Methods [2021]), as well as its capabilities of bridging single-cell data to the corresponding pathological image (Bergensträhle *et al.* [2022], Huang *et al.* [2023]).

With the advancements in ST technologies, there has been a surge of interest in large language models (LLMs) (Devlin *et al.* [2018], Brown *et al.* [2020]). Researchers have recognized the potential of LLMs and have begun leveraging their capabilities to expedite computational biology research, particularly in the field of ST.

For example, Bioinformatics Copilot 1.0, introduced by Wang *et al.* [2024], is a tool powered by a large language model. This tool enables intuitive data analysis through a natural language interface, without requiring programming skills, and supports cross-platform functionality. This tool is with a great potential to accelerate advancements in the biomedical sciences by expediting the data analysis workflow. Also, Choi *et al.* [2024] developed a framework called CELLama that uses a language model to transform cell data into "sentences" that capture gene expressions and metadata, allowing for universal cellular data embedding and analysis. CELLama has the potential to revolutionize cellular analysis by enabling cell typing and analysis of spatial contexts without the need for manual reference data selection or dataset-specific workflows.

There are existing methods that may or may not utilize LLMs as natural language interpreters in the task of decoding ST data, but these methods do involve technologies relevant to LLMs. For example, Luo *et al.* [2024] proposed StereoMM,

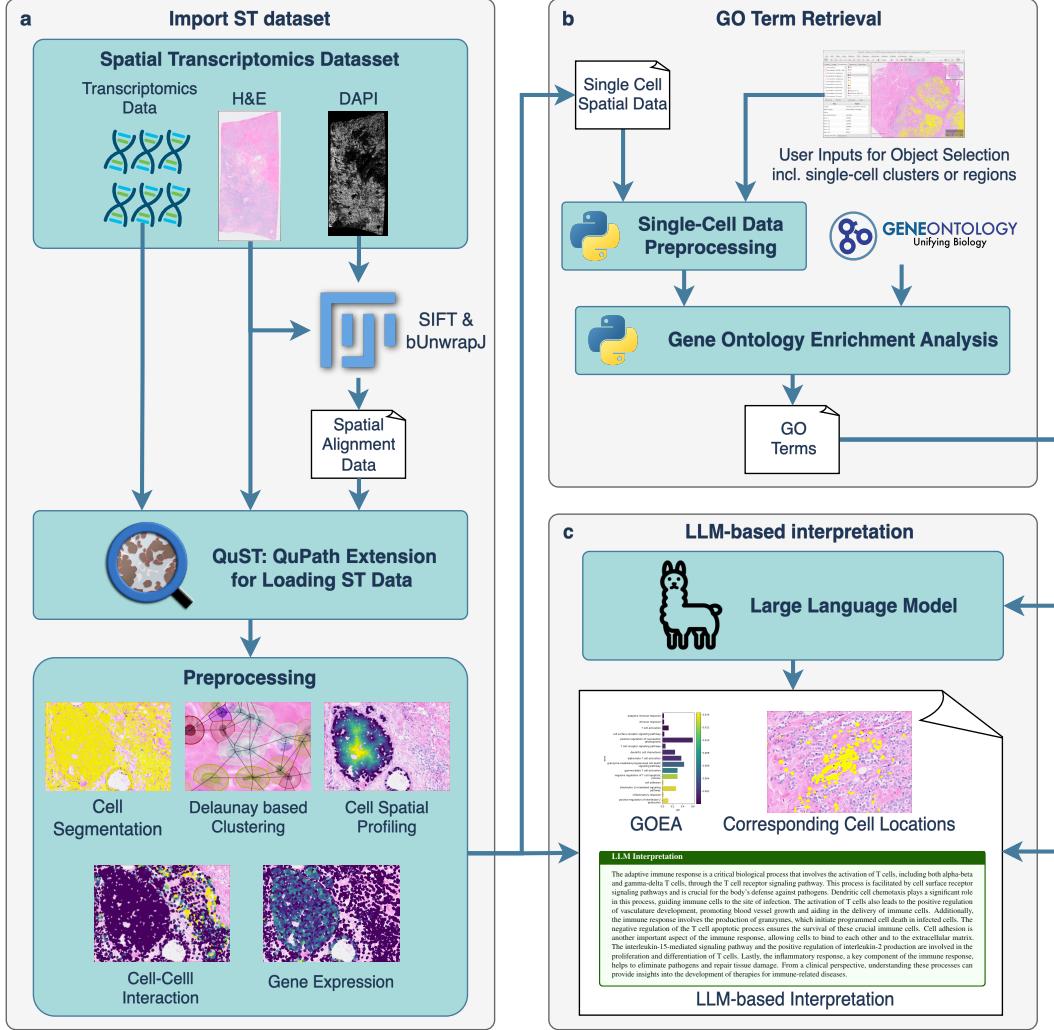


Figure 1: The QuST-LLM workflow for forward analysis includes the following steps: (a), users begin by importing ST data into QuPath using QuST. This step may require additional spatial alignment data, which can be obtained via FIJI if the user is working with a 10x Xenium dataset (see text for more details). Once the ST data is successfully loaded, users can perform analysis and visualization using QuPath and QuST. (b), QuST-LLM takes the objects selected by the user, including single-cell clusters or regions, performs a series of single-cell data preprocessing steps and then obtains a list of GO terms based on GOEA. (c), the spatial data and GO terms are integrated as biological evidence, which can be interpreted using an LLM service. The final outcomes is presented to the users.

a machine learning toolchain that integrates gene expression, histological images, and spatial location data, providing an advanced analysis platform for effectively utilizing multimodal and high-throughput data from spatially resolved omics technologies. Ji *et al.* [2024] proposed SpaCCC, a method for inferring spatially resolved cell-cell communications in ST data using a fine-tuned single-cell language model and a functional gene interaction network. SpaCCC embeds ligand and receptor genes into a unified latent space and identifies likely interacting pairs based on their proximity in this space. Lastly, Cui *et al.* [2024] proposed scGPT, a foundational model for single-cell biology that effectively distills biological insights and can be optimized for superior performance in tasks like cell type annotation, multi-omic integration, and gene network inference.

Despite its potential, the interpretation of ST data presents significant challenges due to the complexity and volume of the information it generates. Advanced computational tools are required to manage, analyze, and interpret these data effectively. Current approaches often rely on a combination of bioinformatics tools to preprocess and analyze the data, but they fall short in providing comprehensive biological insights. To address this need, we present QuST-LLM, an extension of QuPath that enables computational biologists to explore spatial biological problems while providing visualization and analysis capabilities for whole slide image (WSI) analysis.

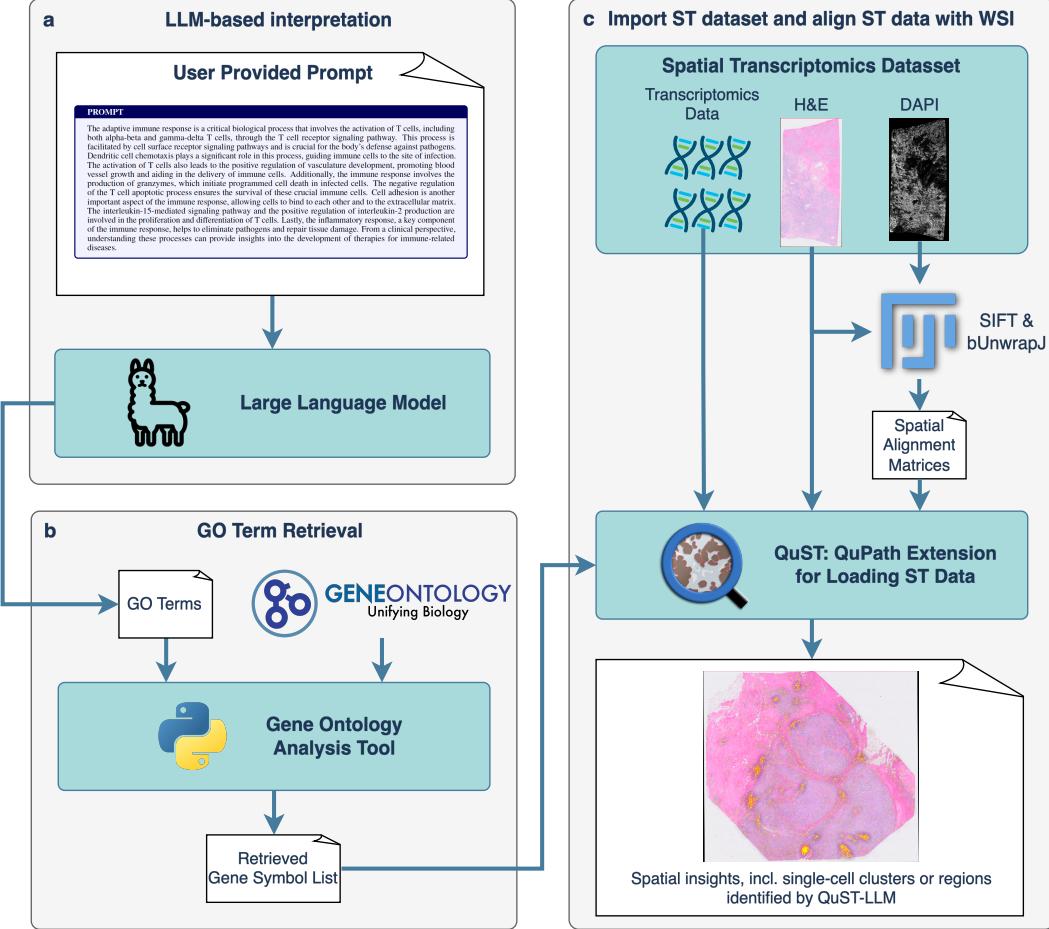


Figure 2: The QuST-LLM workflow for backward analysis includes the following steps: (a), users begin by providing languages describing the required biological evidences. A LLM service is then interpreting the inputs and obtains the key terms which may be used to isolate the sub-graph of the GO. (b), QuST-LLM identifies the key genes by using GOEA based on the obtained GO terms. (c), given the ST data which has been loaded into QuST, the users can then identify the cells which may highly relevant to the sentences provided by the users.

QuST-LLM utilizes a LLM as the backbone. LLMs, such as GPT-4, have demonstrated remarkable abilities in natural language processing, including the interpretation and generation of human-like text based on vast amounts of data. By leveraging LLMs and gene ontology (GO) (Ashburner *et al.* [2000]), a knowledge graph containing , QuST-LLM can translate complex biological annotations into accessible and comprehensive explanations, significantly enhancing the interpretability of ST data.

2 Methods

QuST-LLM relies on two major components: QuPath (Bankhead *et al.* [2017]) and QuST (Huang [2024]). QuPath is an open-source software platform widely used for bioimage analysis, offering powerful tools for visualizing and analyzing high-resolution tissue images. QuST is a powerful extension for QuPath that seamlessly integrates whole slide image (WSI) and ST analysis, providing enhanced capabilities for spatial biology research. It enables the visualization of spatial gene expression data alongside histopathological images, allowing researchers to explore the molecular landscape of tissues at an unprecedented resolution. In this section, we will introduce the analyzing tools and use cases available in QuST.

QuST-LLM facilitates a seamless workflow from data acquisition to biological interpretation. First, users load ST data into QuPath, select regions of interest (ROIs). Next, QuST-LLM takes care of the following analyzing tasks automatically, including profiling gene expression within these regions using SCANPY (Wolf *et al.* [2018]), identifying

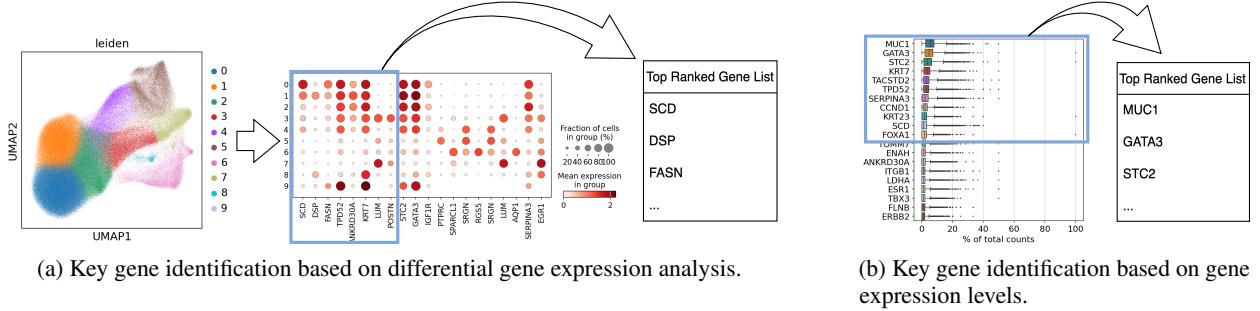


Figure 3: Two approaches for obtaining key genes.

the corresponding genes, and performing GO enrichment analysis (GOEA) using GOATOOLS (Klopfenstein *et al.* [2018]). Finally, the GO terms are interpreted using LLMs, providing detailed explanations of the biological significance of the selected cells.

One of the challenges in analyzing ST data is aligning it with WSI, as different image modalities are involved. QuST addresses this challenge by offering various data loading approaches for different ST data formats, including 10x Visium, 10x Xenium, NanoString CosMX, and more. Each format requires a specific approach for proper alignment.

Once the data is loaded, QuST provides a range of visualization and analysis tools. Researchers can visualize gene expression patterns in specific tissue regions, overlay multiple data layers, and adjust visualization parameters to highlight different aspects of the data. This flexibility enables detailed exploration of the spatial relationships between gene expression and tissue morphology. By leveraging QuST’s capabilities, researchers can gain valuable insights into the spatial biology of tissues, advancing our understanding of cellular organization and function.

There are two scenarios of using QuST-LLM: forward and backward analyses. The two scenarios will be introduced in the following sub-sections.

2.1 Forward Analysis: Interpreting Spatial Data using LLM

In the forward analysis (see Figure 1), the user begins by loading spatial transcriptomics (ST) data into a software tool called QuST. The specific steps involved in this process may vary depending on the chosen data modality. However, in general, the following procedure is typically required:

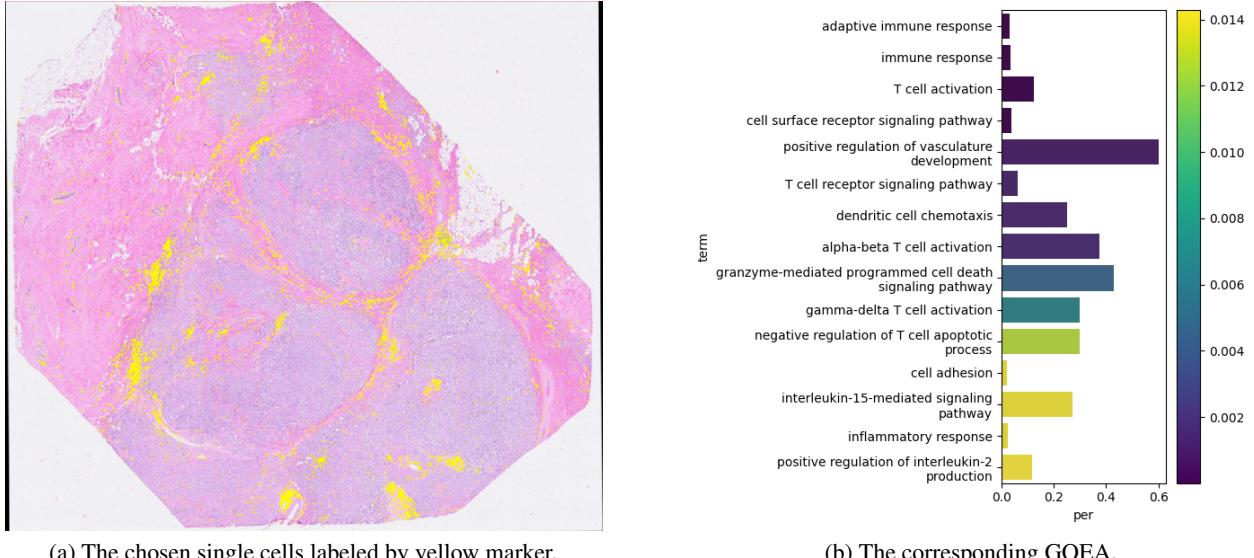
1. Load the high-resolution whole slide image (WSI) into QuPath.
2. Perform necessary analysis on the WSI, including cell segmentation to identify each individual cells.
3. Compute and load spatial correspondence between the WSI and ST data, including the Scale-Invariant Feature Transform (SIFT) (Lowe [2004]) affine matrix and the B-spline basis function coefficients (Sorzano *et al.* [2005]), to QuST.
4. Select the targeting single-cell clusters or regions.
5. Use QuST-LLM to interpret the targeting single-cell clusters or regions.

There are two options available for ST data interpretation using QuST-LLM:

Interpretation based on high ranking key genes. In this option, QuST-LLM interprets the ST data based on key genes which are the genes playing critical roles in the user-specific single-cell clusters or regions (see Figure 3a). In the case if the differential gene expression analysis is insufficient, the gene expression level can be used to measure the ranking of genes (see Figure 3b).

Interpretation based on comparative key genes. In this option, QuST-LLM interprets the ST data based on the differential expression of genes between two or more single-cell clusters or regions defined by the user. This means that the analysis focuses on genes that show significant differences in expression levels between the selected clusters or regions.

These options allow users to gain insights into the spatial distribution and patterns of gene expression within the tissue of interest using the QuST-LLM tool.



LLM Interpretation

The adaptive immune response is a critical biological process that involves the activation of T cells, including both alpha-beta and gamma-delta T cells, through the T cell receptor signaling pathway. This process is facilitated by cell surface receptor signaling pathways and is crucial for the body's defense against pathogens. Dendritic cell chemotaxis plays a significant role in this process, guiding immune cells to the site of infection. The activation of T cells also leads to the positive regulation of vasculature development, promoting blood vessel growth and aiding in the delivery of immune cells. Additionally, the immune response involves the production of granzymes, which initiate programmed cell death in infected cells. The negative regulation of the T cell apoptotic process ensures the survival of these crucial immune cells. Cell adhesion is another important aspect of the immune response, allowing cells to bind to each other and to the extracellular matrix. The interleukin-15-mediated signaling pathway and the positive regulation of interleukin-2 production are involved in the proliferation and differentiation of T cells. Lastly, the inflammatory response, a key component of the immune response, helps to eliminate pathogens and repair tissue damage. From a clinical perspective, understanding these processes can provide insights into the development of therapies for immune-related diseases.

Figure 4: LLM interpretation of high ranking genes based on the selected immuno-cell clusters. (a) The provided whole slide image (WSI) with highlighted single-cell clusters indicated by yellow spots. (b) The results of GOEA, with the x-axis representing the ratio of relevant genes and relevant GO terms, the y-axis showing the list of identified GO terms sorted by uncorrected p-values, and the heat map represents the corresponding p-value for each GO term. (c) The interpretation of the selected immuno-cell clusters as determined by LLM.

2.2 Backward Analysis: Discovering Spatial Insights based on Human Languages using LLM

In the backward analysis (see Figure 2), the user starts from providing a description of the desired targeting single-cell clusters or regions. And then, QuST-LLM will identify these single-cell clusters or regions accordingly. The general procedure is as the following:

1. Load the high-resolution whole slide image (WSI) into QuPath.
2. Perform necessary analysis on the WSI, including cell segmentation to identify each individual cells.
3. Compute and load spatial correspondence between the WSI and ST data, including the SIFT affine matrix and B-spline basis function coefficients, to QuST.
4. User provides a content in human languages describing the required biological evidences.
5. Use QuST-LLM to identify the targeting single-cell clusters or regions.

Algorithm 1 Algorithm to compute the relevance of the given ST data to the user-provided natural language descriptions.

```

1: definition:
2:  $g$ : a string representing a gene symbol defined in GO.
3:  $\mathcal{G} = \{g_1, g_2, \dots\}$ : a set of gene symbols available in the given ST dataset.
4:  $t$ : a string representing a GO term ID.
5:  $x$ : GO category, one of biological process (BP), molecular function (MF) and/or cellular component (CC).
6:  $\mathcal{N}_g^{(x)} = \{t_1, t_2, \dots\}$ : a set as a dictionary that maps  $g$  to a set of associate GO terms based on the given  $x$ .
7:  $\mathcal{M}^{(x)} = \{\mathcal{N}_{g_1}^{(x)}, \mathcal{N}_{g_2}^{(x)}, \dots\}$ : a set of dictionaries  $\mathcal{N}_g^{(x)}$  of the given GO category  $x$ .
8:  $w_g \in \mathbb{R}, \forall g \in \mathcal{G}$ : the weight for  $g$  to be computed,
9:  $\mathcal{W} = (w_{g_1}, w_{g_2}, \dots)$ : a list of  $w_g, \forall g \in \mathcal{G}$ .
10:  $c$ : the symbol representing a cell.

11: input:
12:  $\mathcal{T} = \{t_1, t_2, \dots\}$ : the GO term list obtained via LLM based on the user-provided description in human language.
13:  $\mathcal{C} = \{c_1, c_2, \dots\}$ : list of all targeting cells.
14:  $\mathcal{K}^{(c)} = \{g_1^{(c)}, g_2^{(c)}, \dots, g_k^{(c)}\} \subseteq \mathcal{G}, \forall c \in \mathcal{C}$ : top  $k$  ranked genes for  $c$ .

15: procedure:
16:  $\mathcal{W} \leftarrow \mathbf{0}$  ▷ initial  $\mathcal{W}$ .

17: for each  $x \in \{'BP', 'MF', 'CC'\}$  do ▷ compute  $w_g \in \mathcal{W}, \forall g \in \mathcal{G}$  based on the given  $\mathcal{T}$ .
18:   for each  $g \in \mathcal{G}$  do
19:     if  $\mathcal{N}_g^{(x)} \in \mathcal{M}^{(x)}$  then
20:        $w_g \leftarrow w_g + |\mathcal{N}_g^{(x)} \cap \mathcal{T}|$ 

21:  $\mathcal{W} \leftarrow \mathcal{W} / \|\mathcal{W}\|_1$  ▷ normalize  $\mathcal{W}$ .

22: for each  $c \in \mathcal{C}$  do ▷ compute  $r^{(c)}$  based on  $\mathcal{W}$  and  $\mathcal{K}^{(c)}$ ,  $\forall c \in \mathcal{C}$ .
23:    $r^{(c)} \leftarrow \sum_{g \in \mathcal{K}^{(c)}} w_g$ 

24: output:
25:  $r^{(c)} \in [0, 1], \forall c \in \mathcal{C}$ : the relevance between cell  $c$  and the description provided by the user in human language.

```

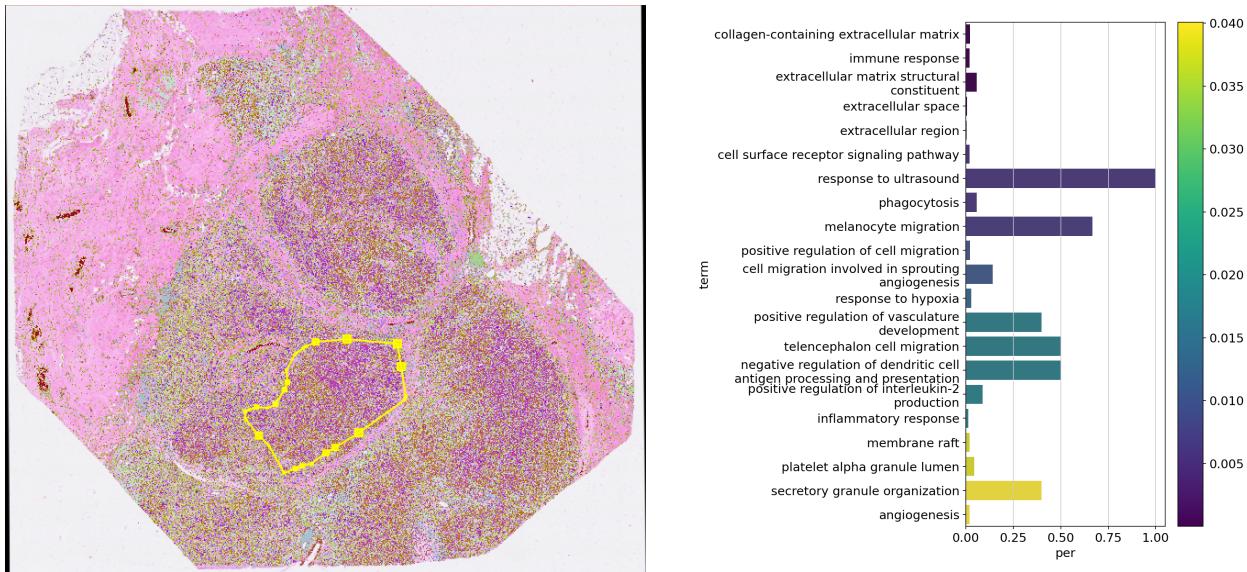
The first three steps are identical to the procedure in the forward analysis as they are essential steps for loading WSI and ST data into QuST. This procedure generates a measurement per cell indicating the level of relevance to the provided description.

In the 4th step, the user selects specific single-cells or regions to target. QuST-LLM performs quality control (QC) and identifies key genes strongly associated with the chosen regions. In the 5th step, QuST-LLM translates the provided description into a list of Gene Ontology (GO) terms linked to a set of genes. Finally, QuST-LLM calculates the correlation between the genes from the identified GO terms and the high ranking genes (either based on their expression levels directly, or based on the correlation levels to some gene sub-groups) in the selected single-cell clusters or regions. The resulting correlation values for each cell are stored in the measurement table in QuPath for further investigation. The pseudo code showing in Algorithm 1 shows how to compute the relevance between a cell and the description provided by the user in human language.

3 Results

In our experiment, we used the data¹ prepared by 10x Genomics that was generated for demonstrating gene expression profiling for formalin-fixed paraffin-embedded (FFPE) human breast samples using the Xenium platform. The sample was obtained as 5µm sections from resected tumor mass tissues of a invasive lobular carcinoma obtained from Avaden Biosciences Inc. The data was generated using 10x Xenium pre-designed panel along with an add-on panel of 100 custom genes. Additionally, the experiment was also conducted using the pre-designed Xenium Human Breast

¹<https://www.10xgenomics.com/datasets/ffpe-human-breast-with-custom-add-on-panel-1-standard>



LLM Interpretation

The biological processes, cellular components, and molecular functions of certain cell groups interact in complex ways to maintain tissue structure and function, with implications for clinical perspectives. The immune response, particularly the adaptive immune response, plays a crucial role in regulating cell migration and monocyte chemotaxis, both of which are essential for tissue repair and inflammation. Calcium-mediated signaling is a key player in these processes, and it also influences the production of interleukin-4, a cytokine involved in immune responses. Interestingly, these processes can be modulated by ultrasound, suggesting potential therapeutic applications.

The extracellular matrix, particularly the collagen-containing component, provides a structural framework for tissues and is involved in various biological processes such as angiogenesis and response to hypoxia. It also plays a role in the organization of secretory granules and the formation of extracellular exosomes, which are involved in cell-cell communication.

The process of angiogenesis, or the formation of new blood vessels, is crucial for tissue repair and regeneration. It involves cell migration and is influenced by factors such as hypoxia and the metabolic process of triglycerides.

Lastly, the migration of melanocytes, cells that produce the pigment melanin, is a key process in skin pigmentation and can be influenced by various factors, including the extracellular matrix and immune responses. Understanding these interactions can provide insights into conditions such as vitiligo, where melanocyte migration is disrupted.

(c) The interpretation generated by LLM.

Figure 5: LLM interpretation of the selected epithelial/tumor-epithelial single-cell clusters based on high ranked gene expression. (a) The provided whole slide image (WSI) with highlighted epithelial/tumor-epithelial single-cell clusters indicated by yellow markers. (b) The result of GOEA, with the x-axis representing the ratio of relevant genes and relevant GO terms, and the y-axis showing the list of identified GO terms sorted by uncorrected p-values. The heat map represents the corresponding p-value for each GO term. (c) The interpretation generated by LLM.

Gene Expression panel. The tissue preparation protocols followed were Xenium In Situ for FFPE - Tissue Preparation Guide and Xenium In Situ for FFPE Tissues – Deparaffinization & Decrosslinking. Post-instrument processing was done according to the Xenium In Situ Gene Expression - Post-Xenium Analyzer H&E Staining protocol.

For annotating regions, due to the fact that the raw H&E image wasn't suitable for more refined exact sub-cellular resolution correspondence due to minute differences in the optics of the microscopy systems, the image was corrected using the methods described at the H&E to Xenium DAPI Image Registration with FIJI Analysis Guide². As a result, the generated SIFT affine matrix and B-spline transformation matrix can be used in the QuST when necessary.

²<https://www.10xgenomics.com/analysis-guides/he-to-xenium-dapi-image-registration-with-fiji>

3.1 Interpreting ST Data using LLM

In this sub-section, we will present the results of LLM-based ST data analysis. The approaches of gene selection include: 1) high gene expression and 2) differential gene expression. The LLM used in this experiment was GPT-4, provided by OpenAI Inc.³.

3.1.1 ST Analysis using LLM-based Approach for High Ranking Gene Expression

Figure 4 illustrates the QuST-LLM interpretation of high ranking gene based on the selected immuno-cell clusters. In the experiment, the genes were selected based on the expression level.

In this experiment, the prompt was:

“Write a paragraph of an integrative and comprehensive summary for the below given key gene ontological terms which are identified by an analysis of a single-cell dataset. Focusing on clinical meaning and structural biology.

{GO_TERM:1, GO_TERM:2, …, GO_TERM:n},”

where the GO terms were obtained using GOATOOOL based on the identified genes of interests in each of the experiments.

In this experiment. QuST-LLM shored detailed insights into the activation of T cells, including alpha-beta and gamma-delta T cells, through the T cell receptor signaling pathway. Additionally, QuST-LLM explores dendritic cell chemotaxis, which guides immune cells to the site of infection, and highlights the role of T cell activation in promoting vasculature development for efficient immune cell delivery. Furthermore, QuST-LLM discusses the involvement of Granzymes in inducing programmed cell death in infected cells, and emphasizes the importance of negative regulation of T cell apoptosis for their survival. Cell adhesion, interleukin signaling, and the inflammatory response are also identified as crucial components of the immune response.

Figure 5 shows the result of en experiment that performed LLM-based interpretation of tumor regions. Given the fact that the chosen regions include various cell types (see Figure 5a), we used differential gene expression analysis to identify the key genes for this experiment.

In this experiment, the prompt was:

“Write a paragraph of a summary for the given content. In the below given content, each row represent some biological processes, cellular components and/or molecular functions of a certain cell group. The summary focuses on the interaction (if any) among these groups. The paragraph is integrative and comprehensive, focus on tissue biological structure and clinical perspectives.

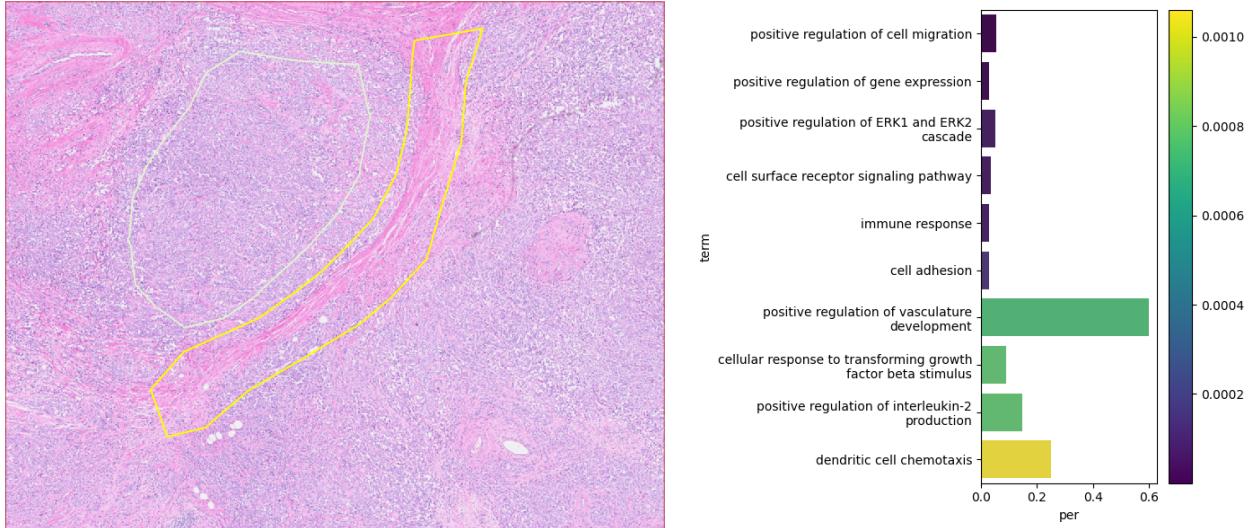
{GO_TERM:1, GO_TERM:2, …, GO_TERM:m},

{GO_TERM:3, GO_TERM:4, …, GO_TERM:n},

…,”

where each row of GO terms represented a category identified using differential gene expression analysis.

The interpretation that QuST-LLM generated suggested the fact that the nature of the chosen regions, including their biological processes, cellular components, and molecular functions interact in complex ways to maintain tissue structure and function, offering potential clinical insights. The the chosen regions, the immune response, especially the adaptive immune response, and calcium-mediated signaling are crucial in regulating cell migration and monocyte chemotaxis for tissue repair and inflammation. Further, the extracellular matrix, particularly its collagen component, provides a structural framework for tissues and contributes to various biological processes including angiogenesis and response to hypoxia. In addition, the extracellular matrix also influences the organization of secretory granules and formation of extracellular exosomes, which are key to cell-cell communication. Understanding the process of angiogenesis and the migration of melanocytes, influenced by factors like the extracellular matrix and immune responses, can provide insights into conditions such as vitiligo.



LLM Interpretation

The positive regulation of cell migration, gene expression, and ERK1 and ERK2 cascade are crucial processes in the biological structure of tissues, playing significant roles in cellular communication, growth, and differentiation. These processes are often mediated by cell surface receptor signaling pathways, which are integral to the immune response and cell adhesion. The positive regulation of vasculature development is another key aspect, contributing to the formation and maintenance of blood vessels within tissues. The cellular response to transforming growth factor beta stimulus is also noteworthy, as it influences various cellular functions, including proliferation, differentiation, and apoptosis. Furthermore, the positive regulation of interleukin-2 production is essential for immune responses, as it promotes the growth and differentiation of T cells, B cells, and natural killer cells. Lastly, dendritic cell chemotaxis, the directed movement of dendritic cells in response to certain chemicals, plays a pivotal role in the immune response, as these cells act as messengers between the innate and adaptive immune systems. From a clinical perspective, understanding these processes can provide insights into the pathogenesis of various diseases and potentially guide the development of therapeutic strategies.

(c) The interpretation generated by LLM.

Figure 6: LLM interpretation of differential gene expression based on the selected tumor-stromal region, in contrast to a tumor region. (a) The provided whole slide image (WSI) with highlighted tumor-stromal region indicated by yellow regions and tumor region indicated by cyan marker. (b) The results of gene ontology enrichment analysis (GOEA), with the x-axis representing the ratio of relevant genes and relevant GO terms, and the y-axis showing the list of identified GO terms sorted by uncorrected p-values. The heat map represents the corresponding p-value for each GO term. (c) The LLM interpretation of the selected tumor-stromal region, in contrast to a tumor region.

3.1.2 ST Analysis using LLM-based Approach for Differential Gene Expression

Figure 6 shows an interpretation based on differential gene expression between a tumor-stromal region and a tumor-epithelial region. QuST-LLM highlighted important processes such as cell migration, gene expression, and the ERK1 and ERK2 cascade in tissue biology. These processes are regulated by cell surface receptor signaling pathways, which are crucial for immune response and cell adhesion. Vasculature development is also a key aspect, contributing to the formation and maintenance of blood vessels. QuST-LLM also provide additional knowledge, *e.g.*, the cellular response to transforming growth factor beta stimulus influences various cellular functions, and the positive regulation of interleukin-2 production is essential for immune responses, dendritic cell chemotaxis, which enables their movement and communication between the innate and adaptive immune systems, further plays a pivotal role, *etc.* Thus, the experiment result suggested that QuST-LLM is able to build a shortcut for understanding these processes, eventually

³<https://openai.com/>

PROMPT

The adaptive immune response is a critical biological process that involves the activation of T cells, including both alpha-beta and gamma-delta T cells, through the T cell receptor signaling pathway. This process is facilitated by cell surface receptor signaling pathways and is crucial for the body's defense against pathogens. Dendritic cell chemotaxis plays a significant role in this process, guiding immune cells to the site of infection. The activation of T cells also leads to the positive regulation of vasculature development, promoting blood vessel growth and aiding in the delivery of immune cells. Additionally, the immune response involves the production of granzymes, which initiate programmed cell death in infected cells. The negative regulation of the T cell apoptotic process ensures the survival of these crucial immune cells. Cell adhesion is another important aspect of the immune response, allowing cells to bind to each other and to the extracellular matrix. The interleukin-15-mediated signaling pathway and the positive regulation of interleukin-2 production are involved in the proliferation and differentiation of T cells. Lastly, the inflammatory response, a key component of the immune response, helps to eliminate pathogens and repair tissue damage. From a clinical perspective, understanding these processes can provide insights into the development of therapies for immune-related diseases.

(a) The prompt used in the experiment.

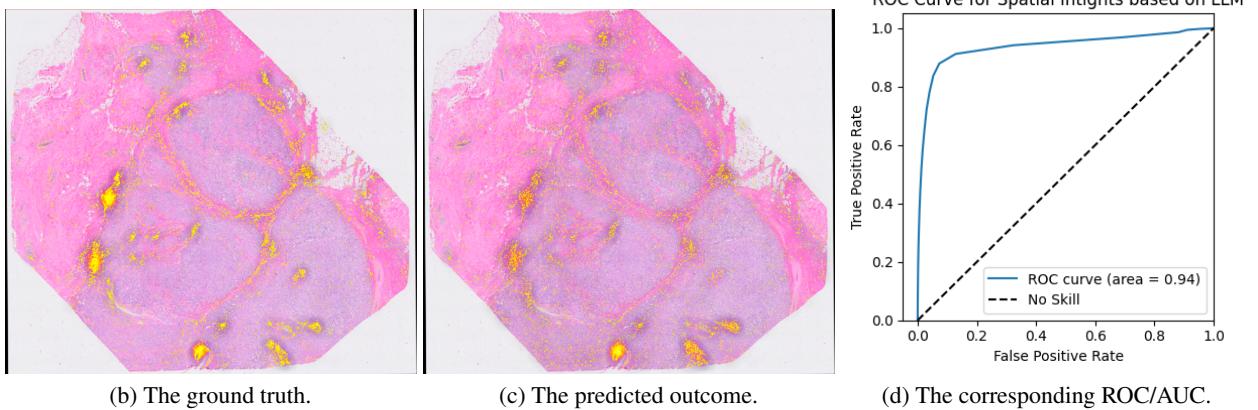


Figure 7: The result showcases the discovery of spatial insights based on human languages using a LLM. (a) is the prompt used in this experiment, which is identical to the LLM-based interpretation showing in Figure 4c. (b) shows the density map of the ground truth indicated by yellow spots. (c) illustrates the density map of the predicted outcome. (d) presents the ROC/AUC curve, indicating the accuracy of the prediction based on the user's description provided in human language.

provides valuable insights into disease pathogenesis and offers potential guidance for the development of therapeutic strategies.

3.2 Discovering Spatial Insights based on Human Languages using LLM

In this sub-section, we will present using LLM to identify single-cell clusters or regions that reflects the description of the given prompt. The prompt used was:

`"Identify corresponding gene ontology term IDs from below content, and return the result in json format.`

`{PROMPT},"`

where the PROMPT was given by the user describing the desired biological status.

Figure 7 demonstrates the experiment using the same prompt as Figure 4c to compare forward and backward analyses outcomes. We identified single-cell clusters via gene expression levels and used the ROC method to gauge model accuracy, yielding a high-quality prediction with an AUC of 0.94, as shown in Figure 7d. The predicted outcome confirmed the role of immune response in T cell activation, vasculature development, and granzyme production, showcasing QuST-LLM's ability to reveal spatial insights using human languages.

4 Conclusion

The integration of large language models (LLMs) into spatial transcriptomics (ST) analysis, embodied in the QuST-LLM tool, signifies a considerable leap in the genomics field. QuST-LLM's ability to transform intricate, high-dimensional ST data into comprehensible biological narratives significantly enhances the interpretability and accessibility of ST data. Our study has demonstrated that QuST-LLM is not only capable of interpreting biological spatial patterns, but it can also identify specific single-cell clusters or regions based on user-provided natural language descriptions. This represents the transformative potential of LLMs in computational biology research and their ability to assist researchers in deciphering the spatial and functional complexities of tissues.

One avenue for future research could be the fine-tuning of the underlying LLM. By training the model on more specific biological and genomic datasets, it could further improve the tool's ability to interpret and analyze ST data. This could lead to more precise interpretations and potentially uncover deeper insights, thereby driving further advancements in biomedical research.

Overall, QuST-LLM represents a significant step forward in making ST data analysis more intuitive and accessible. With the prospect of fine-tuning the LLM, we anticipate even more impactful contributions to the field in the future.

5 Availability

The QuST-LLM is a function provided in QuST, which is developed based on QuPath 0.5.1 and Python 3.10+ and is available under the Apache 2.0 license (<https://github.com/huangch/qust>). A user guide is provided at https://github.com/huangch/qust/user_guide, including a step-by-step tutorials, GPU support, and examples demonstrating the use of the extension in analysis pipelines.

References

- Patrik L. Ståhl, Fredrik Salmén, Sanja Vickovic, Anna Lundmark, José Fernández Navarro, Jens Magnusson, *et al.* Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science*, 353:78–82, 7 2016. ISSN 10959203. doi:10.1126/SCIENCE.AAF2403.
- Amanda Janesick, Robert Shelansky, Andrew D. Gottscho, Florian Wagner, Stephen R. Williams, Morgane Rouault, *et al.* High resolution mapping of the tumor microenvironment using integrated single-cell, spatial and *in situ* analysis. *Nature Communications*, 14, 12 2023. ISSN 2041-1723. doi:10.1038/s41467-023-43458-x.
- Nature Methods. Method of the year 2020: spatially resolved transcriptomics. *Nature Methods*, 18, 1 2021. ISSN 1548-7105. doi:10.1038/s41592-020-01042-x.
- Ludvig Bergenstråhle, Bryan He, Joseph Bergenstråhle, Xesús Abalo, Reza Mirzazadeh, Kim Thrane, *et al.* Super-resolved spatial transcriptomics by deep data fusion. *Nature biotechnology*, 40:476–479, 4 2022. ISSN 1546-1696. doi:10.1038/S41587-021-01075-3.
- Chao-Hui Huang, Yoson Park, Jincheng Pang, and Jadwiga R. Bienkowska. Single-cell gene expression prediction using h&e images based on spatial transcriptomics. volume 12471, pages 17–25. SPIE, 4 2023. ISBN 9781510660472. doi:10.1117/12.2654294.
- Jacob Devlin, Ming Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, 1: 4171–4186, 10 2018.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, *et al.* Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- Yongheng Wang, Weidi Zhang, Siyu Lin, Matthew S. Farruggio, and Aijun Wang. Bioinformatics copilot 1.0: A large language model-powered software for the analysis of transcriptomic data. *bioRxiv*, 2024.04.11.588958, 4 2024. doi:10.1101/2024.04.11.588958.
- Hongyoon Choi, Jeongbin Park, Sumin Kim, Jiwon Kim, Dongjoo Lee, Sungwoo Bae, *et al.* Cellama: Foundation model for single cell and spatial transcriptomics by cell embedding leveraging language model abilities. *bioRxiv*, page 2024.05.08.593094, 5 2024. doi:10.1101/2024.05.08.593094.
- Bingying Luo, Fei Teng, Guo Tang, Weixuan Chen, Chi Qu, Xuanzhu Liu, *et al.* Stereomm: A graph fusion model for integrating spatial transcriptomic data and pathological images. *bioRxiv*, 2024.05.04.592486, 5 2024. doi:10.1101/2024.05.04.592486.

- Boya Ji, Liwen Xu, and Shaoliang Peng. Spaccc: Large language model-based cell-cell communication inference for spatially resolved transcriptomic data. *bioRxiv*, page 2024.02.21.581369, 2 2024. doi:10.1101/2024.02.21.581369.
- Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, *et al.* scgpt: toward building a foundation model for single-cell multi-omics using generative ai. *Nature Methods*, 2 2024. ISSN 1548-7105. doi:10.1038/s41592-024-02201-0.
- Michael Ashburner, Catherine A. Ball, Judith A. Blake, David Botstein, Heather Butler, J. Michael Cherry, *et al.* Gene ontology: tool for the unification of biology. the gene ontology consortium. *Nature genetics*, 25:25–29, 5 2000. ISSN 1061-4036. doi:10.1038/75556.
- Peter Bankhead, Maurice B. Loughrey, José A. Fernández, Yvonne Dombrowski, Darragh G. McArt, Philip D. Dunne, *et al.* Qupath: Open source software for digital pathology image analysis. *Scientific Reports*, 7:16878, 12 2017. ISSN 2045-2322. doi:10.1038/s41598-017-17204-5.
- Chao-Hui Huang. Qust: Qupath extension for integrative whole slide image and spatial transcriptomics analysis. *arXiv:2406.01613 [q-bio.QM]*, 5 2024. URL <https://arxiv.org/abs/2406.01613>.
- F. Alexander Wolf, Philipp Angerer, and Fabian J. Theis. Scanpy: Large-scale single-cell gene expression data analysis. *Genome Biology*, 19:1–5, 2 2018. ISSN 1474760X. doi:10.1186/S13059-017-1382-0/FIGURES/1.
- D. V. Klopfenstein, Liangsheng Zhang, Brent S. Pedersen, Fidel Ramírez, Alex Warwick Vesztrocy, Aurélien Naldi, *et al.* Goatools: A python library for gene ontology analyses. *Scientific Reports 2018* 8:1, 8:1–17, 7 2018. ISSN 2045-2322. doi:10.1038/s41598-018-28948-z.
- David G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60:91–110, 11 2004. ISSN 09205691. doi:10.1023/B:VISI.0000029664.99615.94/METRICS.
- Carlos Ó S. Sorzano, Philippe Thévenaz, and Michael Unser. Elastic registration of biological images using vector-spline regularization. *IEEE Transactions on Biomedical Engineering*, 52:652–663, 4 2005. ISSN 00189294. doi:10.1109/TBME.2005.844030.