# 开放科学浪潮下，图书馆员的新角色

Calvin Wu| 吴非

EBSCO 销售经理

EBSCO Information Services

# 开放科学的过去与现在

**01**

❓ **一项研究的生命周期**

❗ **开放阶段的前移**

# 一项研究的生命周期也许是这样：



Re-use
Ratings
Credits
Citations
Blogs
Tweets

Design

Proposals
Templates
Drafts

Track

Plan

DMPs

Products
Identifiers
Peer Reviews
Versions

Publish,Preserve
Archive

Tracking
Transparency

Collect,Find,
Acquire

Data
Code
Samples
Reagents
Materials
Methods
Instruments
Tools
Subjects

Prepare

Process
Visualize
Analyze

Metadata
Annotations
Formats&Standards
Files
Licenses
Methods&Protocols
Results

Store

Workflow tools
Scripts & Software
Graphics
Model & Simulations

Cloud services
Field Notebooks
ELN
Collaboration spaces

# 一项研究的生命周期也许是这样：

**开放科学的背景：
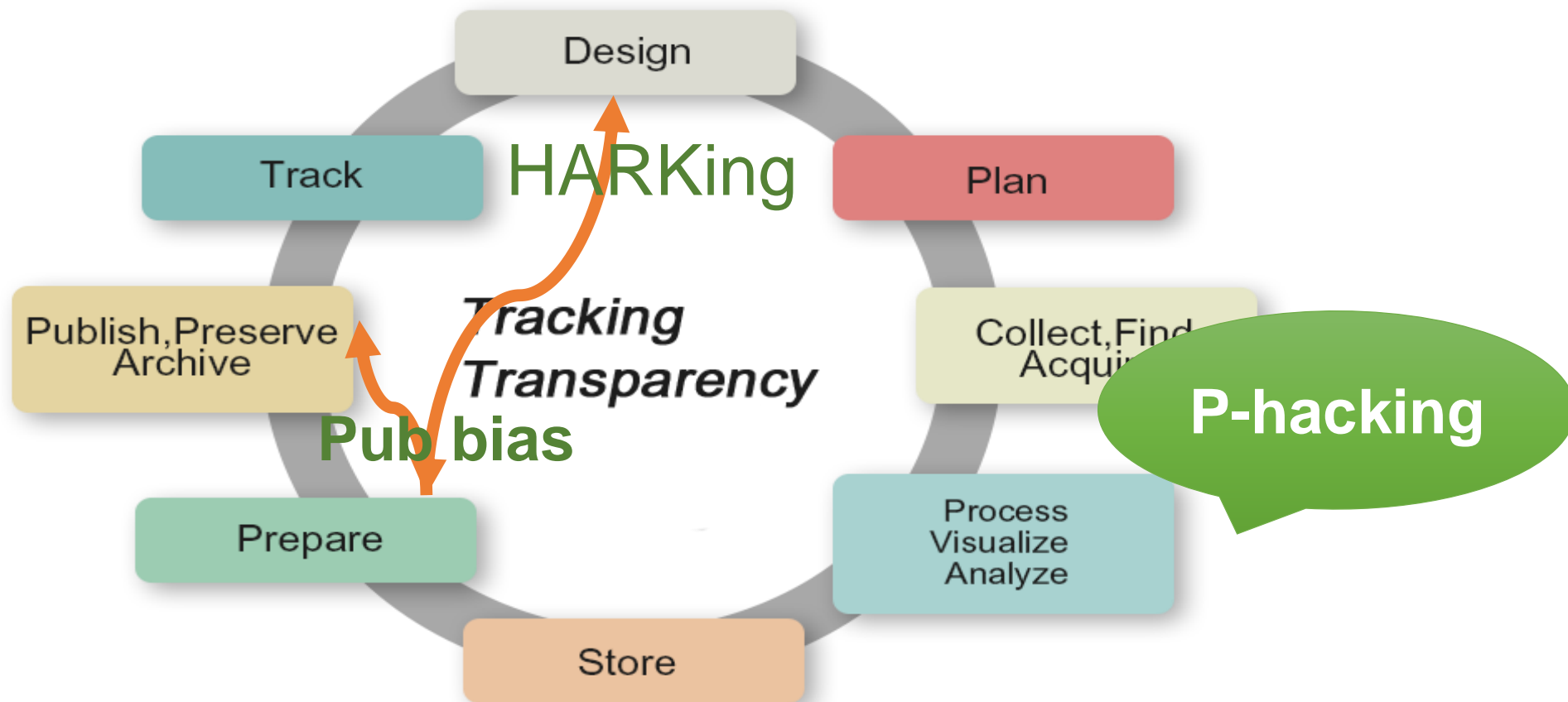顶级期刊纷纷响应，2010s成
为热议话题的"R WORDS"：**

**Repeatability(可重复性)**

**Replicability(可复制性)**

**Reproducibility(可再现性)**

# 如果只分享出版以后的产物，难以避免R-WORDS的危机以及研究不透明的威胁：

今天的开放科学-
期待研究生命周期的每个阶段都会开放，
即便是失败的结果

# 走向开放科学的未来

**02**

- ➢ **开放科学新兴的3D模型**

- ➢ **No raw data, No science, 没有原始数据，就没有科学**

- ➢ **未来图书馆员可能的新角色**

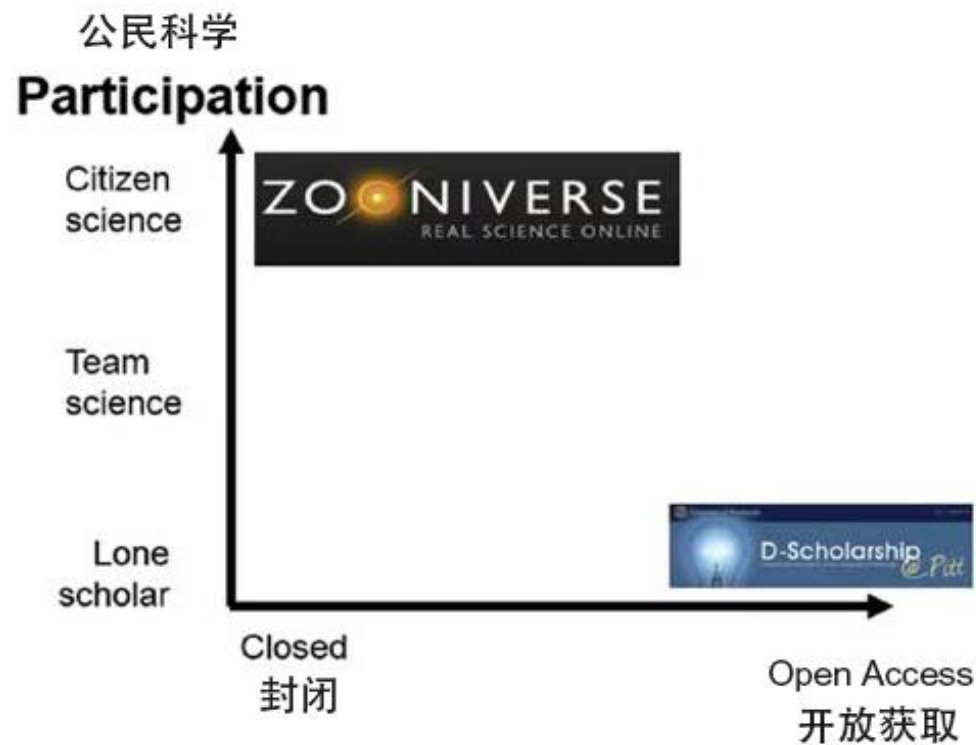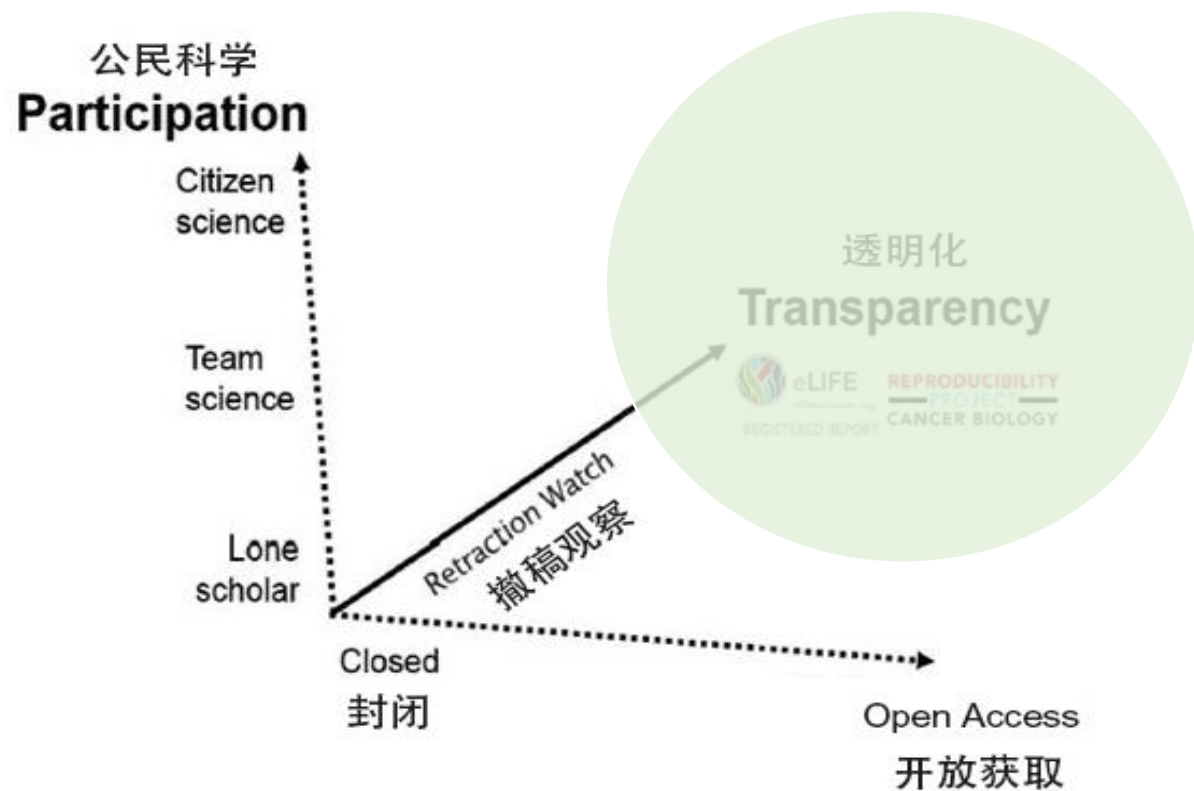Fig. 1: 2D-model of Open Science (Based on the Continuum of Openness in Lyon, 2009).
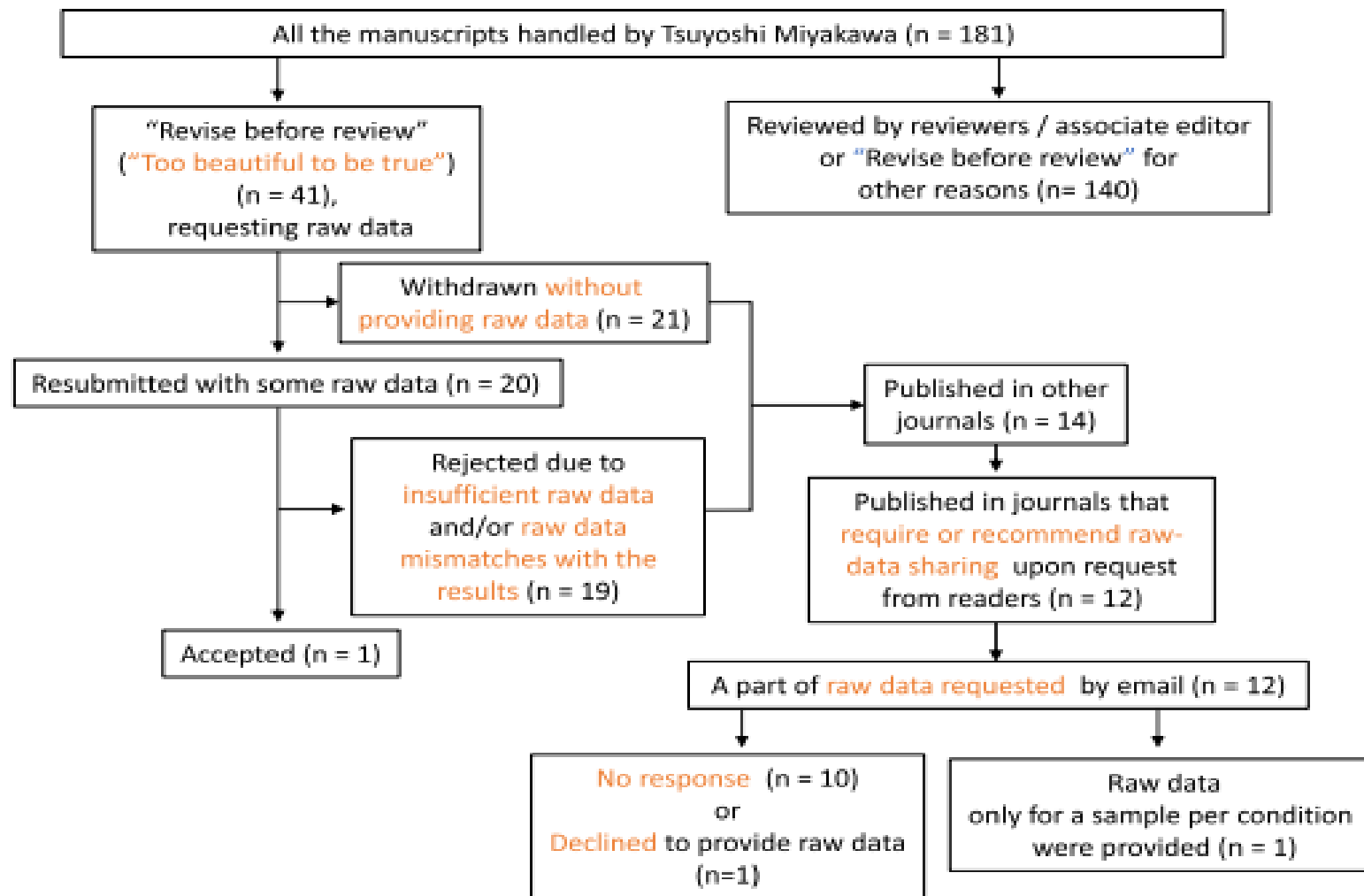
Fig. 2: 3D-Model of Open Science.

# 开放科学新兴的3D模型

**Fig. 1** Flowchart of the manuscripts handled by Tsuyoshi Miyakawa in *Molecular Brain* from December 2017 to September 2019

▌**Tsuyoshi Miyakawa**

原始数据的缺乏或数据捏造是造成不可重复性的另一个可能原因

让原始数据能够公开获取不仅对于数据挖掘、再使用非常重要，而且对于去证实研究的结果是否基于实际的数据，也同样重要

- Tsuyoshi Miyakawa

# 国务院办公厅关于印发科学数据管理办法的通知

「政府预算资金资助形成的科学数据应当按照开放为常态、不开放为例外的原则，由主管部门组织编制科学数据资源目录，有关目录和数据应及时接入国家数据共享交换平台，面向社会和相关部门开放共享。」

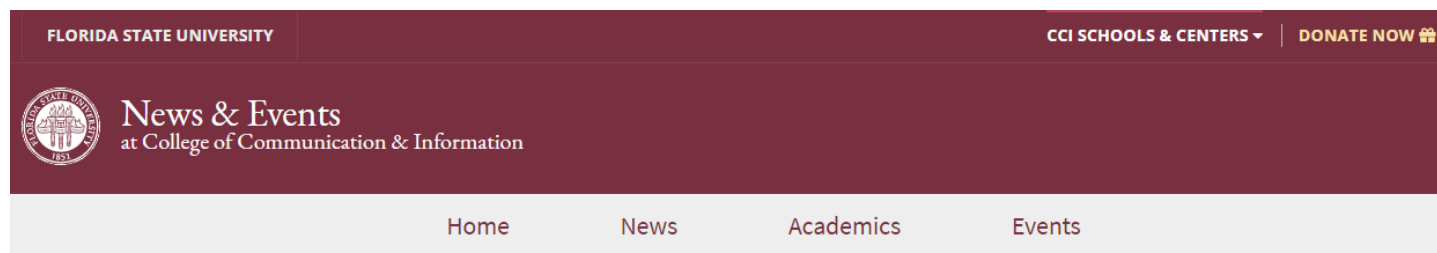http://www.gov.cn/zhengce/content/2018-04/02/content_5279272.htm

# 图书馆员可能的新角色-数据透明化代理人，即Data Librarian@开放科学有关透明化的基本术语

*Table 2: Foundational Terms for Transparency in Open Science.*

| Term | Exposition |
| --- | --- |
| Transparency | The outcome from a suite of behaviours which characterize Reproducible Research |
| Transparency | Facilitates and enhances Research Quality, Research Integrity and Trust |
| Transparency Action | Describes a specific intervention which is a component of the processes, protocols and practices within the Research Lifecycle |
| 数据透明化代理人 | Exemplified by the Data Science roles e.g. 数据图书馆员 . These are key components of the Data Fabric (RDA) and supporting Infrastructure; they promote and demonstrate specific behaviours and practices which lead to culture change towards Open Science |
| Transparency Tool | The software and model frameworks which support Open Science practice |

# 美国学术图书馆出现新型的图书馆职位-研究再现图书馆员

美国佛罗里达大学图书馆2019年9月招聘一位"研究再现"图书馆员



**FLORIDA STATE UNIVERSITY**     CCI SCHOOLS & CENTERS ▾ | DONATE NOW 🎁

News & Events
at College of Communication & Information

Home    News    Academics    Events

## Reproducibility Librarian @ UF Libraries

Posted by Leila Gibradze on September 18, 2019

The George A. Smathers Libraries at the University of Florida seeks a Reproducibility Librarian to develop an institutional strategy for education and support of transdisciplinary research reproducibility and open science. This position will be located in the Health Science Center Library (HSCL) in Gainesville. The Reproducibility Librarian is a year-round tenure-track library faculty position. The person who fills this position leads in designing and implementing a multifaceted program to enhance campus-wide efforts to promote and improve research reproducibility from design to dissemination.  The Reproducibility Librarian participates as an active member of the library- and campus-wide teams to develop programming for and support information retrieval/storage, data science, and research. The incumbent will provide interdisciplinary information consultation services in a variety of modes, design workshops to promote research reproducibility, perform course-integrated instruction, and participate in the Health Science Center Library's teaching program. The position is responsible for special projects as assigned, such service development and evaluation, and development of web-based resources. The librarian works collaboratively in group efforts and maintains professional relationships with faculty, students and colleagues.

**RESPONSIBILITIES**

• Develops a nationally-recognized program in library-based research reproducibility education and support services at the George A. Smathers Libraries in collaboration with colleagues across the libraries

- 在George A. Smathers图书馆与整个图书馆和校园的同事合作，开发一个全国公认的以图书馆为基础的研究可重复性教育和支持服务项目。
- 为佛罗里达大学的学生、教职员工提供研究可重复性和开放科学领域的专业知识和咨询服务。
- 设计并提供关于提高研究可重复性和进行开放科学的技术指导方案，包括实验和计算分析的记录和保存

# 新型图书馆职位-数据服务图书馆员@北卡罗来纳大学格林斯博罗分校2021年2月

## Responsibilities

- Provide outreach and promotion for data services to faculty, students, and the community
- Serve as Libraries' representative on campus-wide data management initiatives and support related policies
- Engage with campus partners to make digital and scholarly data work openly discoverable, accessible, and reusable
- Provide in-depth research consultations on data discovery, data curation, and data management practices
- Provide course-integrated instruction and standalone workshops on data related topics
- Serve as a liaison to relevant disciplines, which includes providing course-integrated information literacy instruction, research consultations, and collection development
- Teach data literacy skills to a variety of audiences and stay abreast of data literacy trends
- Support liaisons in all disciplines with data discovery needs

- 向教师、学生和社区提供数据服务的宣传和推广。
- 作为图书馆的代表参与全校的数据管理活动并支持相关政策
- 与校园伙伴合作，使学术数据可以公开发现、获取和重复使用

ensboro, NC

reensboro has hired for this role

Apply on company website

Save

# 03

## 支持开放科学的工具

# Code Ocean ：

A centralized platform for the creation, sharing, publication, preservation and reuse of executable code and data.

**Code Ocean是一个创建、共享、发布、保存和重复使用可执行的代码和数据的集中平台。**

# Code Ocean ：

Code Ocean 平台可以将研究产出做成一个标准、安全和可执行的研究包，称为「Capsule胶囊」

# **Reproducible Capsule**
# **可再现的胶囊**



Code

Env

Result

Data

# Code Ocean ：



原始程序代码和数据在哪里？

能够让任何人重制研究的确切计算环境是什么？

# Code Ocean :

# Code Ocean :

Journals & Magazines > IEEE Transactions on Signal P... > Volume: 59 Issue: 9 ❓

## Sensing Matrix Optimization for Block-Sparse Decoding

Code Available

**Publisher:** IEEE     [ Cite This ]     [📄 PDF]

Lihi Zelnik-Manor ;  Kevin Rosenblum ;  Yonina C. Eldar    **All Authors**

| 108 Paper Citations | 3 Patent Citations | 2502 Full Text Views |

Abstract

Document Sections

I.   Introduction

II.  Prior Work on Sensing Matrix Design

III. Sensing Matrix Design for Block-Sparse Decoding

IV.  WCM—Weighted Coherence Minimization

**Abstract:**
Recent work has demonstrated that using a carefully designed sensing matrix rather than a random one, can improve the performance of compressed sensing. In particular, a well-designed sensing matrix can reduce the coherence between the atoms of the equivalent dictionary, and as a consequence, reduce the reconstruction error. In some applications, the signals of interest can be well approximated by a union of a small number of subspaces (e.g., face recognition and motion segmentation). This implies the existence of a dictionary which leads to block-sparse representations. In this work, we propose a framework for sensing matrix design that improves the ability of block-sparse approximation techniques to reconstruct and classify signals. This method is based on minimizing a weighted sum of the interblock coherence and the subblock coherence of the equivalent dictionary. Our experiments show that the proposed algorithm significantly improves signal recovery and classification ability of the Block-OMP algorithm compared to sensing matrix optimization methods that do not employ block structure.

# 出版社对于 "Reproducibility可重现性" 的响应



ScienceDirect

Journals & Books

Software Impacts
Open access

Articles & Issues ⌄    About ⌄    Publish ⌄    🔍    ur article ↗    G

R徽章文章代表着可重现，并由CODE OCEAN认证，胶囊会永久性地存放在CODE OCEAN上。

## R-badged Articles

Last update 27 April 2021

This collection presents software publications that are verified for computational reproducibility by the CodeOcean, a cloud-based computational reproducibility platform that helps the community by enabling sharing of code and data as a resource for non-commercial use.
Certified papers have an attached Reproducibility Badge, a permanent Reproducible Capsule and are listed on the CodeOcean website.

Actions for selected articles    Receive an update when the latest issues in this journal are published

# 出版社对于 "Reproducibility可重现性" 的响应

# Code Ocean ：



Allows anyone to find, create and share code and data, without the need for local deployment

允许任何人查找、创建和共享代码和数据，而无需本地部署

# 文章的问题：

**学术文章往往在其研究方法的描述有所欠缺，特别是详尽的步骤、细节或组成的成分**

"The hardest part, by far, was figuring out exactly what the original labs actually did. Scientific papers come with methods sections that theoretically ought to provide recipes for doing the same experiments. But often, those recipes are **incomplete**, **missing important steps**, **details**, or **ingredients**. In some cases, the recipes aren't described at all"

*The* *Atlantic*

# protocols.io

## Many organizations encourage the use of protocols.io

Journals and publishers recommend protocols.io on manuscript submission

AACR — American Association for Cancer Research

PLOS

GSA — Genetics Society of America

eLIFE

**500+ journals**

Funders require or recommend protocols.io in grant guidelines/policies

THE LEONA M. AND HARRY B. HELMSLEY CHARITABLE TRUST

Alex's Lemonade Stand

CZ

GORDON AND BETTY MOORE FOUNDATION

NIH

wellcome

# protocols.io

Allows anyone to find, create and share research methods to support reproducibility and collaborations

允许任何人搜寻、创建和共享研究方法，以支持可重复性与协同合作

# 结语：

# 参考文献：

Miyakawa, T. No raw data, no science: another possible source of the reproducibility crisis. *Mol Brain* 13, 24 (2020). https://doi.org/10.1186/s13041-020-0552-2

LYON, L. Transparency: The Emerging Third Dimension of Open Science and Open Data. Liber Quarterly: The Journal of European Research Libraries, [s. l.], v. 25, n. 4, p. 153–171, 2016. DOI 10.18352/lq.10113. Disponível em: http://search.ebscohost.com/login.aspx?direct=true&db=lls&AN=114854749&site=ehost-live&scope=site. Acesso em: 25 abr. 2021.

Benedikt Fasel (1), Jörg Spörri (2,3), Josef Kröll (3), Kamiar Aminian (1) 2017. Functional calibration for trunk and lower limb fixed inertial sensors. protocols.iohttps://dx.doi.org/10.17504/protocols.io.itrcem6

Butler, D., Karpowitz, C., & Pope, J. (2017). Who Gets the Credit? Legislative Responsiveness and Evaluations of Members, Parties, and the US Congress. *Political Science Research and Methods, 5*(2), 351-366. doi:10.1017/psrm.2015.83

L. Zelnik-Manor, K. Rosenblum and Y. C. Eldar, "Sensing Matrix Optimization for Block-Sparse Decoding," in IEEE Transactions on Signal Processing, vol. 59, no. 9, pp. 4300-4312, Sept. 2011, doi: 10.1109/TSP.2011.2159211.

# THANKS
谢谢