



R语言

数据可视化之美

专业图表绘制指南

张杰 / 著

R语言

数据可视化之美 专业图表绘制指南

张杰 / 著



电子工业出版社
Publishing House of Electronics Industry
北京·BEIJING

内 容 简 介

本书主要介绍使用 R 中的 ggplot2 包及其拓展包绘制专业图表的方法。本书先介绍了 R 语言编程基础知识, 以及使用 dplyr、tidyr、reshape2 等包的数据操作方法; 再对比了 base、lattice 和 ggplot2 包的图形语法。本书首次系统性地介绍了使用 ggplot2 包及其拓展包绘制类别对比型、数据关系型、时间序列型、整体局部型等常见的二维图表的方法, 以及使用 plot3D 包绘制三维图表(包括三维散点图、柱形图和曲面图等)的方法。另外, 本书也首次介绍了论文中学术图表的图表配色、规范格式等相关技能与知识。

未经许可, 不得以任何方式复制或抄袭本书之部分或全部内容。
版权所有, 侵权必究。

图书在版编目(CIP)数据

R 语言数据可视化之美: 专业图表绘制指南 / 张杰著. —北京: 电子工业出版社, 2019.6
ISBN 978-7-121-36366-5

I. ①R… II. ①张… III. ①统计分析—应用软件—指南 IV. ①C819-62

中国版本图书馆 CIP 数据核字(2019)第 070960 号

责任编辑: 石 倩

印 刷: 中国电影出版社印刷厂

装 订: 中国电影出版社印刷厂

出版发行: 电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本: 787×980 1/16 印张: 18.25 字数: 461 千字

版 次: 2019 年 6 月第 1 版

印 次: 2019 年 6 月第 1 次印刷

定 价: 109.00 元

凡所购买电子工业出版社图书有缺损问题, 请向购买书店调换。若书店售缺, 请与本社发行部联系, 联系及邮购电话: (010) 88254888, 88258888。

质量投诉请发邮件至 zlts@phei.com.cn, 盗版侵权举报请发邮件至 dbqq@phei.com.cn。

本书咨询联系方式: 010-51260888-819, faq@phei.com.cn。

前言

本书主要介绍使用 R 中的 `ggplot2` 包及其拓展包绘制专业图表的方法。本书先介绍了 R 语言编程基础知识，以及使用 `dplyr`、`tidyr`、`reshape2` 等包的数据操作方法；再对比了 `base`、`lattice` 和 `ggplot2` 包的图形语法。本书首次系统性地介绍了使用 `ggplot2` 包及其拓展包绘制类别对比型、数据关系型、时间序列型、整体局部型等常见的二维图表的方法，以及使用 `plot3D` 包绘制三维图表（包括三维散点图、柱形图和曲面图等）的方法。另外，本书也首次介绍了论文中学术图表的图表配色、规范格式等相关技能与知识。

本书定位

虽然现在 Python 语言越来越流行，尤其是在机器学习与深度学习等领域，但是 R 语言在数据分析与可视化方面仍然具有绝对的优势，其中 `ggplot2` 包及其拓展包人性化的绘图语法大受用户的喜爱，特别是生物信息与医学研究者。现在 *Nature*、*Science* 和 *Cell* 等期刊上大量的图表都是使用 R 语言绘制的，所以很有必要系统性地介绍 R 语言的绘图方法。

R `ggplot2` 有两本很经典的教程：*ggplot2 Elegant Graphics for Data Analysis* 和 *R Graphics Cookbook*，两书重点介绍了 `ggplot2` 包的绘图语法及常见图表的绘制方法，但是其介绍的图表种类并不多。所以本书基于 R 中的 `ggplot2` 包及其拓展包和 `plot3D` 包，系统性地介绍了几乎所有常见的二维和三维图表的绘制方法，包括简单的柱形图系列、条形图系列、折线图系列，以及复杂的和弦图、矩形树状图、日历图等。

读者对象

本书适用于想学习数据分析与可视化相关专业课程的高校学生，以及对数据分析与可视化感兴

趣的职场人士阅读，尤其是 R 语言用户。从软件掌握程度而言，本书同样适用于零基础学习 R 语言的用户。

阅读指南

全书内容共有 9 章，其中，第 1 章和第 2 章是后面 7 章的基础，第 3~8 章都是独立章节，可以根据实际需求有选择性地进行学习。

第 1 章 介绍 R 语言编程与数据可视化基础，对比了 base、lattice 和 ggplot2 包的图形语法，重点介绍了 ggplot2 包的图形语法；

第 2 章 介绍 R 语言数据处理基础，重点介绍了使用 dplyr、tidyr、reshape2 等包的数据操作方法；

第 3 章 介绍类别比较型图表，包括柱形图系列、条形图系列、南丁格尔玫瑰图、径向柱图等约 30 种图表；

第 4 章 介绍数据关系型图表，包括二维和三维散点图、气泡图、等高线图、三维曲面图、三相相图、二维和三维瀑布图、相关系数热力图等约 60 种图表；

第 5 章 介绍数据分布型图表，包括一维、二维和三维的统计直方图和核密度估计图、抖动散点图、点阵图、箱形图、小提琴图等约 50 种图表；

第 6 章 介绍时间序列型图表，包括折线图和面积图系列、日历图、螺旋图系列、量化波形图、地平线图约 20 种图表；

第 7 章 介绍局部整体型图表，包括饼图、散点复合饼图系列、旭日图、矩形树状图、马赛克图、华夫饼图等约 20 种图表；

第 8 章 介绍高维数据的可视化方法，包括分面图系列、矩阵散点图、热力图、平行坐标系图、RadViz 图、图标法等约 20 种图表；

第 9 章 介绍论文中学术图表的常用技能，包括常见的截图与图片处理软件及其功能、矢量图片的修改、论文中学术图表数据的提取与重绘、论文中学术图表的规范与调整等。

应用范围

本书的图表绘制方法都是基于 R 中的 ggplot2 包及其拓展包和其他绘图包实现的，几乎适应于所有常见的二维和三维图表。但是由于依据《地图管理条例》第十五条规定：“国家实行地图审核制度。向社会公开的地图，应当报送有审核权的测绘地理信息行政主管部门审核。但是，景区图、街区图、地铁线路图等内容简单的地图除外。”本书本来有专门的章节讲解使用 R 语言如何绘制不同地理坐标

投影下，从世界到不同国家与区域的地图，但是由于地图审核周期等方面的原因忍痛移除。

适用版本


本书所用 R 版本为：3.3.3。R 作为开源免费的软件，数据分析与可视化的包更新迭代很快，这是它的优势。但是有时候有些代码运行可能会由于 R 或者 R 包版本的更新，而出现函数弃用（deprecated）的情况。此时，需要自己更新代码，使用新的函数替代原有的函数等。


源代码

本书配备有几乎所有图表的 R 语言源文件及其.csv 或.txt 格式的数据源文件。但是需要注意的是，如果运行的 R 语言版本没有安装相应的数据分析与可视化的包（package），那么请预先安装相应的包，才能成功运行代码。同时，也请注意运行 R 语言和 R 包的版本是否已经更新。


与我联系

因本人知识与能力所限，书中纰漏之处在所难免，欢迎并恳请读者朋友们给予批评与指正，可以通过邮箱联系笔者本人；如果读者有关于 R 语言学术图表或商业图表绘制的问题，可以联系笔者交流。另外，更多关于 R 语言图表绘制的教程请关注笔者的博客、专栏和微博平台。也可以重点关注我们的微信公众号：EasyCharts，也可以添加笔者微信：EasyCharts。R 语言数据分析与可视化的文章会优先发表在微信公众号平台。

 邮箱：easycharts@qq.com

 知乎专栏：<https://zhuankan.zhihu.com/EasyCharts-R>（知乎账号：EasyCharts）

 博客：<http://easychart.github.io/>

 新浪微博：https://weibo.com/easycharts?source=blog&is_all=1（微博账号：EasyCharts）

致谢

桃李春风一杯酒，江湖夜雨十年灯。笔者的处女作《Excel 数据之美：科学图表与商业图表的绘制》也至今出版逾两年，一直想着要修订这本书的。但是旧书未翻新，新书忙于码字改稿，实在是愧于读者。其实，在撰写这本新书的时候，数次想放弃。写书实在是一件费力劳神的事情，笔者是凭借着对数据可视化的热爱才坚持至今。

这本书从 2017 年 5 月 25 日开始动笔，断断续续居然也花费了两年的时间。与其说是花费，不如说是陪伴吧。笔者经常对朋友开玩笑说，心情不好的时候码码代码、画画图表，是一件消磨时间、

放松心情的事情。

在断断续续的写稿中，笔者也认识了很多热爱数据分析与可视化的朋友，甚是荣幸，也得益于他们的帮助。很感谢《R 语言游戏数据分析与挖掘》的作者谢佳标老师和先锋信息科技有限公司 CEO 林祯舜老师对笔者的鼓励与帮助，也因此有幸参加 2018 年的 R 语言大会；也非常感谢在码字、写代码时一起交流学习的李誉辉（四川大学高分子学院）、杜雨（美团用户平台—大数据与算法部—商业分组部）、刘钰（河南大学土木建筑学院）、厚缦（深圳中观经济咨询有限公司）等诸多技术大佬。因为有你们的帮助，所以才有今天这本书。

最后，想对大家说，也是对自己说：且将新火试新茶，诗酒趁年华！

作者

2019 年 3 月 31 日

读者服务

轻松注册成为博文视点社区用户（www.broadview.com.cn），扫码直达本书页面。

- **下载资源**：本书如提供示例代码及资源文件，均可在 [下载资源](#) 处下载。
- **提交勘误**：您对书中内容的修改意见可在 [提交勘误](#) 处提交，若被采纳，将获赠博文视点社区积分（在您购买电子书时，积分可用来抵扣相应金额）。
- **交流互动**：在页面下方 [读者评论](#) 处留下您的疑问或观点，与我们和其他读者一同学习交流。

页面入口：<http://www.broadview.com.cn/36366>





目 录

第 1 章 R 语言编程与绘图基础	1
1.1 学术图表的基本概念	2
1.1.1 学术图表的基本作用	3
1.1.2 学术图表的基本类别	5
1.1.3 学术图表的绘制原则	7
1.2 你为什么要选择 R	8
1.3 R 软件的安装与使用	15
1.3.1 R 与 RStudio 的安装	15
1.3.2 包的安装与加载	16
1.4 R 语言编程基础	17
1.4.1 数据类型	17
1.4.2 数据结构	18
1.4.3 数据属性	21
1.4.4 数据的导入导出	23
1.4.5 控制语句与函数编写	26
1.5 R 语言绘图基础	27
1.6 ggplot2 图形语法	29
1.6.1 geom_×××()与 stat_×××()	31
1.6.2 视觉通道映射	34

1.6.3	度量调整	37
1.6.4	坐标系	43
1.6.5	图例	52
1.6.6	主题系统	54
1.6.7	位置调整	57
1.7	学术图表的色彩运用原理	61
1.7.1	颜色模式	61
1.7.2	颜色主题的搭配原理	66
1.7.3	学术图表的颜色主题	69
1.7.4	颜色方案的拾取使用	71
1.7.5	颜色主题的应用案例	74
1.8	图表的基本类型	77
1.8.1	类别比较	78
1.8.2	数据关系	79
1.8.3	数据分布	79
1.8.4	时间序列	80
1.8.5	局部整体	80
1.8.6	地理空间	81
第 2 章 R 语言数据处理基础		83
2.1	表格的转换	84
2.1.1	表格的变换	84
2.1.2	变量的变换	85
2.1.3	表格的排序	86
2.2	表格的整理	86
2.2.1	表格的拼接	86
2.2.2	表格的融合	87
2.2.3	表格的分组操作	89
第 3 章 类别比较型图表		92
3.1	柱形图系列	93
3.1.1	单数据系列柱形图	94

3.1.2 多数据系列柱形图	95
3.1.3 堆积柱形图	96
3.1.4 百分比堆积柱形图	97
3.2 条形图系列	98
3.3 不等宽柱形图	99
3.4 克利夫兰点图系列	100
3.5 坡度图	102
3.6 南丁格尔玫瑰图	103
3.7 径向柱形图	107
3.8 雷达图	109
3.9 词云	112
第4章 数据关系型图表	116
4.1 散点图系列	117
4.1.1 趋势显示的二维散点图	117
4.1.2 分布显示的二维散点图	124
4.1.3 气泡图	128
4.1.4 三维散点图	130
4.2 曲面拟合图	133
4.3 等高线图	136
4.4 切面图	138
4.5 三元相图	139
4.6 散点曲线图系列	141
4.7 瀑布图	143
4.8 相关系数图	149
4.9 韦恩图	151
4.10 树形图	152
4.11 圆堆积图	154

4.12	和弦图	156
4.13	桑基图	160
第5章	数据分布型图表	163
5.1	统计直方图和核密度估计图	165
5.1.1	统计直方图	165
5.1.2	核密度估计图	165
5.2	数据分布系列	169
5.2.1	散点分布图系列	170
5.2.2	柱形分布图系列	172
5.2.3	箱形图系列	173
5.2.4	其他图表	178
5.3	二维统计直方图和二维核密度估计图	188
5.3.1	二维统计直方图	188
5.3.2	二维核密度估计图	188
5.4	金字塔图和镜面图	192
第6章	时间序列型图表	194
6.1	折线图与面积图系列	195
6.1.1	折线图	195
6.1.2	面积图	195
6.2	日历图	199
6.3	螺旋图	202
6.4	量化波形图	207
6.5	地平线图	210
第7章	局部整体型图表	213
7.1	饼状图系列	214
7.1.1	饼图	214
7.1.2	圆环图	216
7.1.3	复合饼图系列	216

7.2	旭日图	219
7.3	矩形树状图	221
7.4	马赛克图	224
7.5	华夫饼图	227
第 8 章	高维数据可视化	229
8.1	高维数据的变换展示	231
8.1.1	主成分分析法	231
8.1.2	t-SNE 算法	233
8.2	分面图	234
8.3	矩阵散点图	238
8.4	热力图	240
8.5	平行坐标系图	243
8.6	RadViz 图	245
8.7	图标法	246
8.7.1	基于星形图的图标法	247
8.7.2	基于柱形图的图标法	249
8.7.3	切尔诺夫脸谱图	251
8.8	表格图	254
第 9 章	论文中学术图表的升级技能	255
9.1	图片的截取与处理软件	256
9.1.1	常见截图软件	256
9.1.2	图片处理软件	256
9.2	论文中学术图表的规范与调整	257
9.2.1	图片的格式与转换	260
9.2.2	图片的分辨率	262
9.2.3	图片的色彩要求	264
9.2.4	图片的物理尺寸	265
9.2.5	图片的标注格式	266

9.2.6 图片的占内存容量	266
9.2.7 在 R 中导出图表	268
9.3 图表绘制的必备技能	269
9.3.1 矢量图表元素的修改	269
9.3.2 期刊论文的图片提取	271
9.3.3 图表数据的重新提取	271

参考文献	274
------------	-----



第1章

R 语言编程与绘图基础

1.1 学术图表的基本概念

学术图表是为论文结论 (conclusion) 提供证据的视觉方式。所以, 论文作者为了产生强烈的视觉效果, 应该通过分析实验数据, 精心设计可视化图表。本书开篇先跟大家讲讲学术图表的类型。通常学术论文中主要有三类图表, 如图 1-1-1 所示。流程示意图和数据展示图都是非常讲究技能的图表, 本书重点讲解的是数据展示图。

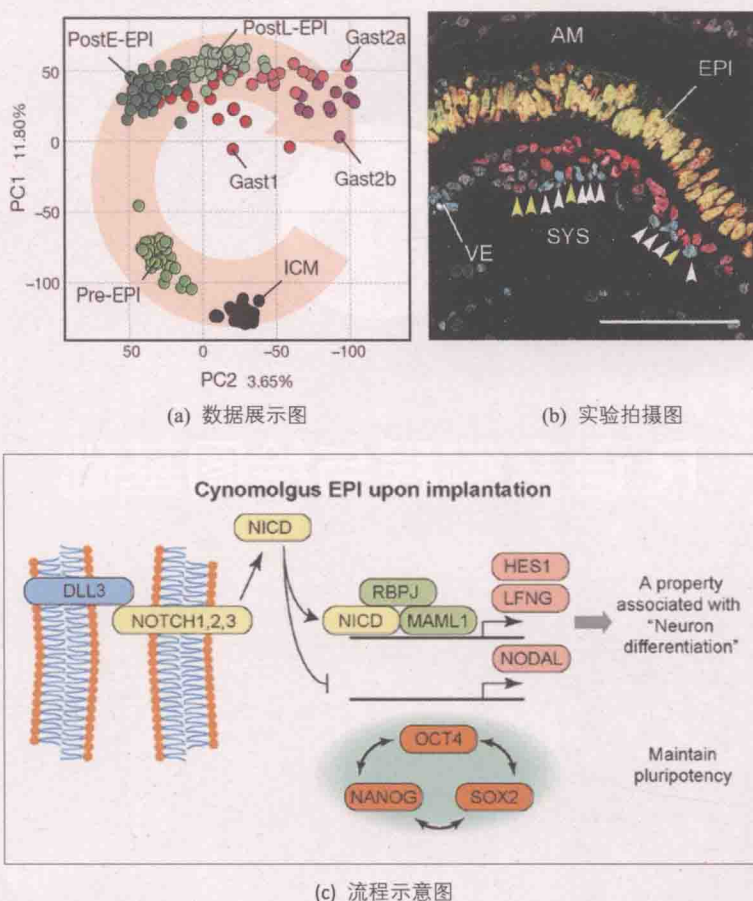


图 1-1-1 不同类型的图表^[1]

1. 数据展示图: 先根据数据绘制成图表, 再将其导出生成图片, 主要包括各种点线图、柱形图、饼图等统计图表, 一般使用 Excel、GraphPad Prism、SigmaPlot、Origin、MATLAB、Python、R 等专业绘图软件绘制 (Excel 并非如大众所说不能导出高分辨率的图片和矢量图)。注意, 保存图片时,

一定要保存成高分辨率的 TIFF 格式和 EPS 矢量格式的图片，因为矢量图片是可以使用图片处理软件进行再编辑的。由数据生成的图表是可重复修改的，因此一定要保存好原始数据，一旦发现图表有任何问题就可以进行修改。

2. 实验拍摄图：使用设备或者仪器拍摄采集的图片，包括显微镜、扫描仪及摄像机等所拍照片。一定要在最刚开始时就拍成高清的（设置成高分辨率），也就是要保证原始图片的高分辨率，接下来处理图片就会比较方便，免得因为图片质量不佳而重复实验。必要的话，把每张图片存储成 TIFF 和 JPG 两种格式（以应对部分期刊的特殊要求）。

3. 流程示意图：使用简明的线条、基本图形和箭头等绘制论文中的重要实验流程或步骤，用以说明基本原理或解释文字材料，一般使用 PPT、Visio、Illustrator、CorelDRAW、3DMax 等软件绘制。

1.1.1 学术图表的基本作用

图表在学术论文中是很重要的一部分。实验结果通常是论文的核心和主要部分，而实验结果一般以图表的形式呈现。读者经常通过图表来判断这篇文章是否值得阅读，所以每个图表都应该能不依赖正文而独立存在。所谓“一图抵千言”（A picture is worth a thousand words）。图表设计是否精确且合理直接影响数据的完整与准确表达，从而影响论文的质量。图表是期刊评审过程中仅次于摘要的关键一环，准确而美观的图表能促进审稿人和读者对论文表达的快速理解。以 *Nature* 上的文章 *Cotranslational signal-independent SRP preloading during membrane targeting*^[2] 选取的前两页为例（见图 1-1-2），我们首先关注的是论文的标题（title），其次是第一页最开始的摘要（abstract），接下来我们就被这些包含大量实验数据与信息的图表所吸引。在每页的文章中，包含图名（figure）的图表部分几乎占据整个页面的 1/4~1/3，由此可见图表在论文中的重要性。

根据 Edward R. Tufte 的 *The Visual Display of Quantitative Information*^[3] 和 *Visual Explanations*^[4] 的阐述，图表在论文的作用主要有：

- （1）真实、准确、全面地展示数据；
- （2）以较小的空间承载较多的信息；
- （3）揭示数据的本质、关系、规律。

第三点作用尤为重要，Matthew O. Ward 也提出，可视化的终极目标是洞悉蕴含在数据中的现象和规律，这包括多重含义：发现、决策、解释、分析、探索和学习^[5]。表 1-1-1 所示的原始数据是 31 组 x - y 的二维数据。仅仅只从数据的角度去观察数据，就很难发现 x 与 y 之间的具体关系。将实际的数据分布情况使用二维可视化的方法呈现，如图 1-1-3 所示，则可以快速地从数据中发现数据内在的模式与规律。所以，有时使用数据可视化的方法也可以很好地帮助我们去分析数据。

LETTER

doi:10.1038/nature14106

Cotranslational signal-independent SRP preloading during membrane targeting

Anton W. Chamone¹, Katherine C. L. Hahn¹ & Jadhav Tyagi^{1,2}

Ribosome-associated factors must precisely divide the limited information available to nascent polypeptides to direct them to their correct cellular fate. It is unclear how the low-complexity information encoded by the nascent chain suffices for such a diverse range of cellular fates. We investigated the cotranslational membrane-targeting cycle using ribosome profiling in yeast cells coupled with biochemical fractionation of ribosome populations. We show that the SRP preferentially binds nascent ERGA before their targeting signals are translated. Non-coding mRNA elements can pre-load the signal-independent pre-encounter of SRP over which the complex kinetically transitions between degradation in the cytosol and attachment to the polypeptide and mRNAs that mediate SRP-substrate selection and membrane targeting.

Secretory proteins are proposed to target to the endoplasmic reticulum (ER) membrane either co- or post-translational by signal sequence transduction^{1–3}. Mechanistic models of ER targeting and the role of the SRP derive primarily from cell-free systems using model proteins^{4–6}, using the quantification of how flow-pulse-labeled proteins in the ER⁷. To investigate membrane-targeting in vivo, we used ribosome-associated factors to map the cotranslational membrane-targeting cycle using ribosome profiling in yeast cells coupled with biochemical fractionation of ribosome populations (Fig. 1a). We detected a cotranslational membrane-targeting cycle in which SRP preferentially binds nascent ERGA before their targeting signals are translated. Non-coding mRNA elements can pre-load the signal-independent pre-encounter of SRP over which the complex kinetically transitions between degradation in the cytosol and attachment to the polypeptide and mRNAs that mediate SRP-substrate selection and membrane targeting.

Secretory proteins are proposed to target to the endoplasmic reticulum (ER) membrane either co- or post-translational by signal sequence transduction^{1–3}. Mechanistic models of ER targeting and the role of the SRP derive primarily from cell-free systems using model proteins^{4–6}, using the quantification of how flow-pulse-labeled proteins in the ER⁷. To investigate membrane-targeting in vivo, we used ribosome-associated factors to map the cotranslational membrane-targeting cycle using ribosome profiling in yeast cells coupled with biochemical fractionation of ribosome populations (Fig. 1a). We detected a cotranslational membrane-targeting cycle in which SRP preferentially binds nascent ERGA before their targeting signals are translated. Non-coding mRNA elements can pre-load the signal-independent pre-encounter of SRP over which the complex kinetically transitions between degradation in the cytosol and attachment to the polypeptide and mRNAs that mediate SRP-substrate selection and membrane targeting.

of translation on soluble versus membrane-bound ribosomes. For cytosolic proteins, soluble ribosome-protected reads were distributed across the entire reading frame, consistent with complete translation in the cytosol (Extended Data Fig. 4). For secretory proteins, both soluble and membrane-bound ribosomes produced protected reads. Cytosolic transcripts represented only a small fraction of genes secretory proteins, apart from the secretory mRNA pool was membrane-associated. By the classic understanding of cotranslational targeting, secretory proteins should be the membrane only after exposing a targeting signal⁸. Thus, there should be fewer ERGA bound on the membrane translating the protein of transcripts not set

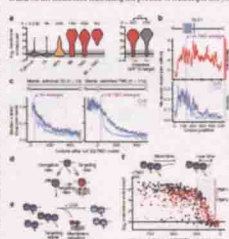


Figure 1 | Cotranslational membrane-associated ERGA. (a) Schematic of the cotranslational membrane-targeting cycle. (b) Ribosome profiling reads for ERGA and non-coding RNA. (c) Kinetic transition of SRP-ERGA complex.

targeted, that is, at codon positions upstream of the first SS or TMD. However, the membrane-bound ribosomes protected reads were preferentially distributed across the entire reading frame (Fig. 1b and Extended Data Fig. 4). This suggests that once targeted, secretory mRNAs remain associated to the ER and their translation continues at the membrane. This is consistent with the observed presence of secretory ERGA on the membrane before synthesis of the targeting signal⁸. The small fraction of secretory mRNA in cytoplasmic pools targeted ERGA is probably represents the protein round of targeting.

The presence of membrane-bound mRNA proteins may also be due to low secretory transcripts on targeted to the membrane. The higher read density for these mRNAs suggests that the signal sequence was exposed by the ribosome, as expected from cotranslational signal-dependent targeting of soluble ERGA to the membrane (Fig. 5b) and Extended Data Fig. 4d). Surprisingly, the bound reads after signal exposure was gradual, extending to non-coding RNA that contained soluble for hundreds of codons after SS or TMD exposure. This result was inconsistent with the degenerate attractor bias errors proposed for the SRP^{9,10} and suggests that degradation upstream on cytosolic ERGA upon exposure of a targeting signal (the Suppressor Decoy).

The idea that there is a kinetic competition between cotranslational targeting to the cytosol and ERGA targeting to the membrane has two viable predictions (Fig. 5b). First, pharmacological inhibition of degradation with cycloheximide (CHX) should disrupt these processes, enhancing targeting to the membrane component ERGA and promoting their depletion from the soluble fraction (Fig. 5b). Cells were subjected to a brief, non-toxic CHX incubation before ERGA read analysis of soluble and membrane-bound ribosomes. Importantly, such brief incubation did not perturb nascent secretory proteins (Extended Data Fig. 4). By contrast, CHX treatment markedly reduced the soluble secretory reads, but only after cytosolic ERGA exposure the first SS or TMD that is consistent with the model (Fig. 5b and Extended Data Fig. 4b).

The kinetic competition between targeting and degradation predicts that cotranslational membrane attachment is influenced by translation termination. In the absence of a degradation arrest, the probability of ERGA reading the membrane correspondingly will increase as the first SS or TMD is bound, either at the terminus (Fig. 5b) or Extended Data Fig. 4b). Indeed, we observed a decrease in membrane-bound membrane enrichment of secretory ERGA when the first targeting signal is near the C-terminus. Thus, cotranslational membrane enrichment of SS or TMD must be targeted to the SRP posttranslational (Supplementary Discussion). Overall, our data suggest that cotranslational targeting to the ER is not an unimpeded or a slower round of translation as soluble ribosomes that establishes a pool of ERGA binding mRNAs that mature translation at the membrane (Extended Data Fig. 4).

We next determined which ERGA are substrates of the SRP in vivo. Immunoprecipitation of ERGA from soluble ERGA was followed by ribosome profiling of both SRP-associated polypeptides and monosomes (Fig. 5c and Extended Data Fig. 2a). Few transcripts encoding cytosolic or mitochondrial proteins were enriched on SRP, limiting to specific highly ERGA-associated transcripts. The SRP bound to all secretory ERGA that were cotranslationally targeted to the membrane, including SRP-dependent and SRP-independent proteins (Fig. 2b, c).

The number of ribosome-protected reads from soluble, SRP-bound transcripts diminished after ribosome exposure of the first SS or TMD, as expected from their targeting function (Fig. 2a). The use of gradual and rapid SRP-ERGA cytosolic inhibition will alter the targeting signal but not after exposure to the cytosol. This supports the notion that degradation proceeds on cytosolic ribosomes over an SRP binds to contrast with the expected SRP-induced degradation arrest. Indeed, blocking degradation with CHX for 2 min before translation marked degradation in SRP-bound reads, but only for ERGA

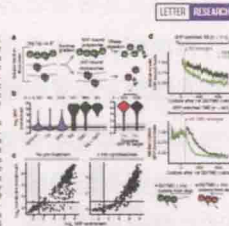


Figure 2 | Cotranslational membrane-associated ERGA. (a) Ribosome profiling reads for ERGA and non-coding RNA. (b) Kinetic transition of SRP-ERGA complex. (c) Kinetic transition of SRP-ERGA complex.

exposing their first targeting signal (Fig. 2b). In principle, the altered targeting of soluble ERGA to the membrane after SRP-induced degradation could result in a delay in SRP binding rather than a lack of degradation arrest. Comparing SRP and membrane enrichment to transcripts indicated that this is not the case. ERGA encoding targeting signals that do not target to the ER membrane (Supplementary Discussion, Fig. 2b and Extended Data Fig. 2b). Addition of CHX allowed the two signal ERGA to reach the membrane, including the SRP-ERGA complex as compared to ERGA targeting. We conclude that the SRP binds the nascent chain quickly and cotranslational degradation causes termination of ERGA before targeting.

Although degradation arrest is not a general consequence of SRP binding in vivo, we next used allowed that a rare quality-directed avoidance of degradation facilitates SRP binding^{9,10}. An intrinsic, non-SRP-dependent degradation avoidance should increase ribosome-protected reads at the same codon in both soluble SRP-bound and membrane-bound polypeptides. Indeed, several transcripts presented such local increases in ribosome-protected reads at sites corresponding to exposure of a targeting signal on the ribosome (Extended Data Fig. 3a–c). Distinct kinetostatic attachment mechanisms observed at these sites included 'beams of rap codons'¹¹ and stable hydrophobic domains, such as stretches of positively charged amino acids, or proline motifs, particularly within the SS region. While these secretory transcripts were not significantly enriched in these domains compared to the non-secretory ERGA (Fig. 3d), the few non-secretory proteins that cotranslationally targeted to the SRP were enriched in degenerate attractor elements present at sites that exposed a rare signal hydrophobic sequence for SRP binding (Extended Data Fig. 3d, 4). We speculate that the presence of such elements enhances SRP recognition of the rare signal hydrophobic motifs to these non-secretory proteins. To understand the basis for the specificity of the SRP in vivo, we next determined the initial point of SRP recruitment to ribosomes

(a) (b)

图 1-1-2 论文摘取的页面案例¹¹

表 1-1-1 四组二维数据点集 (相同的 x 变量, 不同的 y 变量: y1, y2, y3, y4)

x	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
y1	4.6	5.4	5.2	6.6	5.9	6.1	5.8	6.8	6.5	6.7	6.9	11.1	8.2	10.3	12.8	13
y2	6.1	11.6	16.6	19	22.7	31.8	34	33.7	35.6	34.5	39.6	58.3	57.7	72.9	68.4	82.6
y3	5.5	31.1	33.1	51.8	55.7	60.7	63.5	75.5	84.4	84.6	76.3	92.4	81.6	91	88.1	93.8
y4	1	3	4.9	7.9	9.8	12	18.9	24.7	28.9	28.6	39.3	33.2	42.1	54.4	43.3	90.2

x	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	-
y1	20.8	12.4	15.9	15.3	38.8	35.9	24.3	54.5	62.9	43.8	76.9	91	96.9	51.4	100	-
y2	84.5	82	89.1	102.1	68.1	96.3	108.5	76.7	107.6	103.4	116.5	106.4	142.5	115.1	110.5	-
y3	101.3	103	107.4	104.3	110.7	103.4	113.6	105.1	112.5	119.3	113.7	109.5	108.7	110.1	118.8	-
y4	81.2	90.8	70.9	66.8	67.5	88.6	116.9	141.4	104	161.4	101.8	137.1	175.3	119.5	257.3	-