

委托协议编号

## 技术服务（测试化验加工）委托协议

委托任务名称 1100 例人血浆样本蛋白质组学检测

甲方 北京市心肺血管疾病研究所

单位通讯地址 北京市朝阳区安贞路 2 号

乙方 上海中科新生命生物科技有限公司

单位负责人 陈薇

联系人 蒋春丽

联系电话 18301047312

单位通讯地址 上海市闵行区园美路 58 号 1 号楼 15 楼

签订日期：2025 年 6 月 12 日

签订地点：北京

有效期限：2025 年 6 月 31 日 至 2026 年 6 月 31 日

## 填写说明

- 一、本协议适用于我院科研人员在项目研究过程中支付给外单位的检验、测试、化验及加工等费用时需要签署的协议。
- 二、合同封面的委托任务名称指本合同的测试加工等具体内容，应用简明规范的专业术语明确概括所要完成的服务内容。
- 三、本合同的甲方和乙方名称，须按单位公章的详细名称填写，若涉及外文名称，首次出现时应写明全称及简称。
- 四、本协议书未尽事项，可由当事人附页另行约定，并可作为本协议的组成部分。如协议研究内容涉及国家秘密或重大商业秘密的，双方应另行签署保密义务。
- 五、使用本协议书时约定无须填写的条款，应在该条款处注明“无”等字样。
- 六、协议书要求 A4 纸打印，一式 5 份，左侧装订，正文内容所用字型应不小于 5 号字，协议正本中所涉及与本协议约定事项有关的技术资料及其指定附件备齐后应合装成册，其规格大小应与协议书一致。
- 七、乙方需提供测试化验加工的原始数据，甲方务必保留原始数据 10 年以上以备审计抽查。
- 八、协议需法人或委托代理人签署意见后加盖医院公章方可生效。

依据《中华人民共和国民法典》及本协议书相关的科研项目、经费管理办法规定，为完成甲方承担的研究任务，经双方协商一致，各方在真实、充分地表达各自意愿的基础上，就本协议书中所描述的委托内容、经费支付、保密内容、知识产权等问题达成如下协议，签订本合同并由签约双方共同恪守。

## **第一条 委托工作的主要内容、加工方式和要求**

### **1、测试加工内容**

甲方委托乙方就 1100 例人血浆样本进行 DIA 全息扫描蛋白质组学相对定量检测

1.1 本项目实验分析流程包括以下步骤：

1.1.1 人血浆 DIA 相对定量检测：DDA 分析采用纳升流速 HPLC 系统 Easy-nLC 1200 进行色谱分离。样品进样到 C18 色谱分析柱(Thermo Scientific, ES802, 1.9 $\mu$ m, 75 $\mu$ m\*20 cm) 进行线性梯度分离（0.1%甲酸乙腈水溶液（乙腈为 84%）），流速为 300 nL/min。纳升级高效液相色谱分离后的样品用最新一代的高分辨质谱 Orbitrap Astral 进行 DDA 质谱分析。

1.1.2 人血浆 DIA 相对定量检测：采用 Thermo Scientific Vanquish Neo UHPLC+Orbitrap Astral 24min(最新一代高通量质谱，实现血液蛋白质组超深度鉴定)的全息扫描蛋白质组学科研服务。

1.1.3 组学检测：本项目分析流程主要包括 DDA 建库与 DIA 分析两个阶段。质谱实验分析流程主要包括蛋白质提取、肽段酶解、色谱分级、液相色谱-串联质谱（LC-MS/MS）DDA 数据采集、数据库检索等步骤；正式实验阶段主要包括 DIA 分析，质控分析，定性定量结果分析及生物信息学分析。下机数据采用 Spectronaut 进行定性定量分析，提供蛋白定性定量列表，肽段定性定量列表，GO, KEGG, PPI 等生物信息学分析结果。

### **2、测试加工方式和要求**

2.1 技术服务的方式：甲方提供样本，乙方完成全部检测工作。

2.2 技术服务的要求：客观检测，符合数据的质量要求。

## **第二条 考核指标及验收方式**

双方确定以下列标准和方式对乙方的技术服务工作成果进行验收：

1. 乙方完成技术服务工作的形式：按照合同要求客观检测。

2. 技术服务工作成果的验收标准：达到合同中技术服务质量要求。
3. 技术服务工作成果的验收方法：成果报告以纸质版和电子版两种形式发送给甲方。经甲方签字确认后，验收报告生效。
4. 验收地点：北京市心肺血管疾病研究所

第三条 测试化验加工细目：

序号	测试化验加工的内容	测试结果的呈现方式	计量单位	单价 (万元/单位)	数量	金额(万元)
1	1100 例人血浆样本 DIA 蛋白质组学相对定量检测	电子版和纸质版的成果报告	例	0.1	1100	110
合计						110

#### 第四条 经费支付方式:

1. 委托应支付费用共计 110 万元, 由甲方提供。
2. 支付方式一次: (一次或分期) 支付乙方 (按以下第 ③种方式):

①一次总付: \_\_\_\_\_万元。乙方在甲方付款前, 即需提供测试服务。

②分期支付:

第一次支付\_\_\_\_\_万元, 甲方在合同签订后\_\_\_\_日内支付。

第二次支付\_\_\_\_\_万元, 甲方在乙方全部测试技术服务完成并通过验收后\_\_\_\_日内支付。

③其它方式:

合同签订完成 30 日内, 乙方向甲方提供 2 份保函, 其中: 合同总价 5% (伍万伍仟元整人民币) 的履约保函, 保函期限为一年, 全部服务完成经甲方验收合格后退还; 另外合同总价 5% (伍万伍仟元整人民币) 的质量保函, 保函期限为两年, 待验收签字确认合格之日起免费售后服务执行 12 个月后 (若售后服务无问题) 退还。

甲方收到乙方开户银行保函后甲方向财政办理合同支付手续。甲方支付费用 7 日前, 乙方应将对应金额的法定发票提供甲方审核, 待审核通过后甲方按照合同约定向乙方支付费用, 如发票审核不合格, 或者乙方未按规定提供保函的, 甲方有权延期支付费用。

#### 第五条 知识产权归属

1. 双方在申请本课题之前各自所获得的知识产权及相应权益均归各自所有, 不因共同申请本课题而改变。
2. 本协议所产生的所有成果的知识产权全部归属于甲方, 乙方不得利用测试结果单独申报任何形式的成果。
3. 在课题执行过程中各自向对方提供的相关信息, 不构成向对方授予任何关于知识产权的许可行为。
4. 本合作协议不在各方之间建立任何商业上的代理、合作关系。

#### 第六条 保密条款

1. 乙方保证不向甲方以外的人员提供或披露本合同的委托内容及未公开的信息和资料。包括但不限于本协议的委托内容及结果。
2. 双方保证采取一切合理和必要措施和方式对委托中知悉的对方商业秘密进行保密。

---

## 第七条 承诺

1. 如委托的任务涉及人类遗传资源采集、收集、买卖、出口、出境等，乙方承诺遵照《人类遗传资源管理暂行办法》相关规定执行。
2. 如委托任务涉及动物实验，乙方承诺自觉遵守《实验动物管理条例》，严格选用符合要求的合格动物进行实验，保障动物福利。
3. 如委托任务的研究对象涉及人类受试者，乙方承诺在签署协议前已经将委托任务的实施方案呈交单位伦理委员会讨论，并获得了伦理委员会批准。甲方在完成委托任务的过程中，自觉遵守国内外相关的医学伦理准则，保障保护受试者的安全和权益。
4. 在乙方从事委托事项中发生的不可归责于甲方的人身、财产损害，由乙方自行承担。
5. 乙方保证与甲方无直接经济利益关系，并保证委托关系及事项真实有效。

## 第八条 不可抗力

1. 本协议所指不可抗力是指不能预见、不能避免并不能克服的客观情况，包括但不限于地震、火灾、水灾、战争、政府行为等。
2. 乙方因不可抗力不能履行协议的，应当在不可抗力事件发生之日起七日内将不可抗力事由以书面方式通知甲方，并应当在合理期限内提供证明。
3. 因不可抗力不能履行本协议的，根据不可抗力的影响，部分或全部免除责任。乙方延迟履行后发生不可抗力的，不能免除责任。

## 第九条 违约责任

1. 如无正当理由，甲方未能按期拨付工作经费，且经乙方催促仍不能拨付或不能给出合理解释的，乙方有权暂停履行受托任务。如甲方违约行为给乙方造成损失的，甲方还应承担相应赔偿责任。
2. 如乙方在完成委托工作时出现弄虚作假情况、不履行本协议或履行义务不符合要求的，甲方有权追回全部已拨经费。如乙方违约行为造成甲方损失的，甲方有权要求赔偿并追究乙方相关责任人员的法律责任。
3. 非因甲方违约或非因不可抗力，乙方不能完成受托任务或乙方逾期不能提交全部产出成果的，甲方有权解除本委托。委托解除后，乙方应返还甲方已经拨付的项目经费。如乙方的违约行为给甲方造成损失的，乙方还应承担相应的赔偿责任。

- 
4. 乙方在合作期间及合作结束后，未经甲方书面同意，不得在任何形式的宣传材料、广告、媒体发布或公开声明中，使用甲方的名称（包括全称和简称等）、商标、标志、域名、产品或服务进行宣传或暗示其与甲方存在任何形式的合作关系，包括但不限于技术合作、业务往来、信用担保等。如违反本条内容，甲方有权要求乙方停止此侵权行为，并要求乙方赔偿甲方由此遭受的损失（包括直接损失及间接损失）。

#### **第十条 协议的变更、终止及解除**

1. 本协议的变更应由双方协商一致后达成变更协议，并作为本协议的附件。
2. 本协议可由双方协商一致予以终止。

**第十一条 争议解决：**如在履行本协议的过程中发生争执，双方当事人应友好协商解决，如协商不成，任何一方可向甲方签署地（甲方所在地）有管辖权的人民法院提起诉讼。

#### **第十二条 其他约定事项（如无其他事项，请填“无”）**

根据甲方研究需要，可免费提供 20 例样本的检测。

**第十三条** 本协议一式五份，甲方四份，乙方一份，具有同等法律效力。

**第十四条** 本合同及附件、招标文件，投标文件，中标通知书为本合同不可分割的一部分，与本合同具有同等法律效力。

与本协议约定事项有关的技术资料附件清单：见附件

第十四条 签字盖章页

委 托 方 ( 甲 方)	单位名称	北京市心肺血管疾病研究所 (盖章)		
	法定代表人 或授权代表	蔡军 (签字)		
	经办人	李冰洁 (签字)	经办人 联系电话	13552427617
乙 方	单位名称	上海中科新生命生物科技有限公司 (盖章)		
	法定代表人 或授权代表	陈薇 (签字)		
	经办人	蒋春丽 (签字)	经办人 联系电话	18301047312
	开户名称	上海中科新生命生物科技有限公司		
	开户银行	农行漕河泾开发区支行		
	银行账号	03390800040007818		

附件一 投标分项报价表

序号	服务内容	单价 (元)	数量	合价 (元)
1	人血浆样本 DIA 蛋白质组学 相对定量检测	1000.00	1100	1100000.00
总价				1100000.00

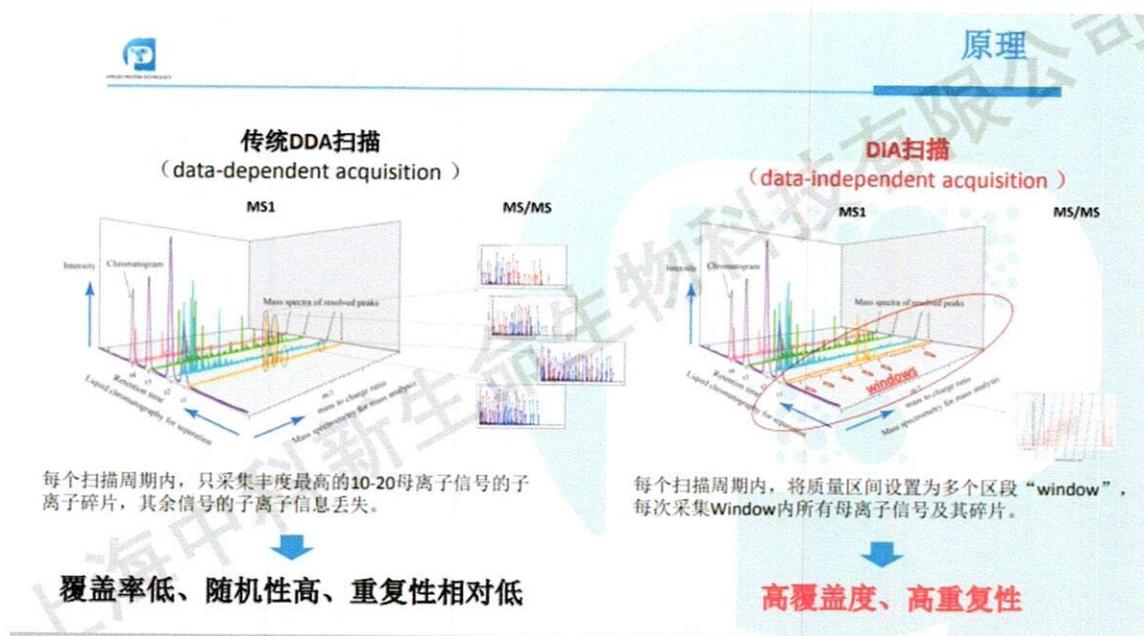
## 附件二 技术方案

### 1 DIA 蛋白组学技术原理

#### 1.1 DIA 蛋白组学技术原理

蛋白组学是通过液质联用技术对蛋白质进行质谱定量分析的。该技术通过大规模分析肽段所产生的质谱数据，比较不同样品中相应肽段的定量信息，从而对肽段对应的蛋白质进行相对定量。DIA (data-independent acquisition) 技术是近年来发展起来的一种新的质谱技术。与传统的 DDA (data-dependent acquisition) 质谱技术相比，DIA 采用了不同的数据扫描模式：将质谱整个全扫描范围分为若干个窗口，然后对每个窗口中的所有离子进行检测、碎裂，从而无遗漏、无差异地获得样本中所有离子的信息。

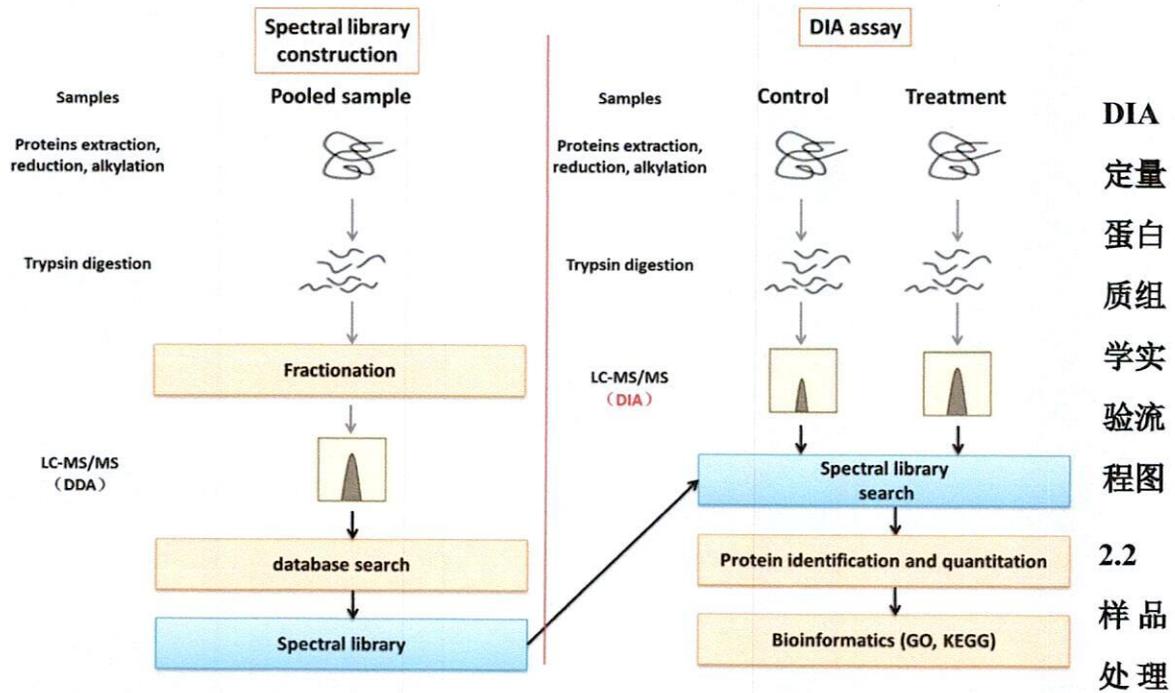
与 DDA 技术相比，DIA 技术的优势包括：（1）采集所有的离子信息，实现更高的数据覆盖度；（2）减少采集的随机性，实现极高的检测重现性、稳定性；（3）采用碎片离子定量，定量精密度、准确性、线性范围大大提高。基于上述技术优势，DIA 技术尤其适用于大规模样本的高度覆盖、稳定和可追溯地分析。



### 2 DIA 全息扫描蛋白质组学整体方案

#### 2.1 分析流程

本项目分析流程主要包括 DDA 建库与 DIA 分析两个阶段。质谱实验分析流程主要包括蛋白质提取、肽段酶解、色谱分级、液相色谱-串联质谱 (LC-MS/MS) DDA 数据采集、数据库检索等步骤对血浆样本分析; 正式实验阶段主要包括 DIA 分析, 质控分析, 定性定量结果分析及生物信息学分析。检测流程示意图如下:



## 和预实验评估及分析方法

### 2.2.1 实验材料准备

#### 2.2.1.1 样本采集的一般性原则

**【一致性原则】**: 每例样本取样的部位、方式、预处理方法需要保持一致

**【快速原则】**: 请务必提前设计准备好实验和材料, 快速取样和分装。

**【分装原则】**: 为避免样本反复冻融, 建议样本采集后立即进行分装。

**【联合分析原则】**: 尽可能保证样本同一批次, 组数及生物学重复一致, 并对不同组学样本进行分装。

**【低温原则】**: 在采集过程中, 请冰上操作, 分离好的样本液氮速冻, 取出保存于 $-80^{\circ}\text{C}$ 冰箱中。

#### 2.2.1.2 样本的包装和运输指南

[1] 样品尽可能采用 1.5ml 或者 2ml 离心管（进口离心管）保存，运输时采用封口膜密封离心管（如管内为有机溶剂，务必采用螺旋口的冻存管并密封）。离心管上标记清楚样品名称后，按顺序整齐排列在冻存盒中。将冻存盒中样品存放的顺序信息对应填写《APT 科研项目送样表》（电子版）。

[2] 不方便存储在离心管中的体积较大的组织样品，推荐采用锡箔纸等材料仔细包装，标记清楚样品名称，按照组别整理整齐，放置在密封袋中。

[3] 推荐采用双层泡沫盒密封包装，盒中加入足量的干冰。

### 2.2.2 样品处理和预实验评估

客户自行收集样品，低温下运输至技术中心。取 4  $\mu$ L 磁珠加入 100  $\mu$ L 血浆中，磁架上孵育 2h。除去废液，用洗涤缓冲液洗涤 5 次。每例样品取适量蛋白，混合成 Pool 样品，用于构建 Spectral Library。所有样品，包括混合 Pool 样品按照 APT 内部 SOP 用 Trypsin 进行溶液内酶解。加 DTT 至终浓度 20mM，30° 反应 2h，冷却至室温，加入 适量 IAA 至终浓度为 25mM，600rpm 振荡 1min，避光室温 30min，加入适量  $\text{NH}_4\text{HCO}_3$  buffer (50mM)将 UA 浓度稀释至低于 1.5M。加入 40  $\mu$ L  $\text{NH}_4\text{HCO}_3$  buffer (2  $\mu$ g Lys-C)，600rpm 振荡 1min，37°C 4h，然后往样品中加入 2  $\mu$ g Trypsin，37°C 16h。脱盐冻干后用 0.1%FA 复溶。OD280 测定肽段浓度。取低丰度肽段，采用 HPRP 方法进行分级，收集所有组分。每个组分肽段冻干后，用 10  $\mu$ L 0.1%FA 复溶，并用 OD280 测定肽段浓度。然后分别取出 2  $\mu$ g 肽段，掺入适量 iRT 标准肽段，进行 DDA 质谱检测。

### 2.2.3 建库方法

DDA 分析采用纳升流速 HPLC 系统 Easy-nLC 11100 进行色谱分离。样品进样到 C18 色谱分析柱(Thermo Scientific, ES802, 1.9 $\mu$ m, 75 $\mu$ m\*20 cm)进行线性梯度分离(0.1%甲酸乙腈水溶液(乙腈为 84%))，流速为 300 nL/min。纳升级高效液相色谱分离后的样品用最新一代的高分辨质谱 Orbitrap Astral 进行 DDA 质谱分析。检测模式：正离子。分 20 级。每级时长 90min，一级质谱扫描范围：350-1650 m/z，质谱分辨率：60,000 (@m/z 1100)，AGC target: 3e6，Maximum IT: 50 ms，动态排除时间：10s。每次一级 MS 扫描(full MS scan)后根据 inclusion list 采集 20 个 ddMS2 扫描(MS2 scans)。

### 2.2.4 DIA 分析方法

---

DIA 分析，每个样本时长 24min。每例样品分别掺入适量 iRT 标准肽段，每个样品进行 1 次 DIA 质谱测试。DIA 分析采用 Thermo Scientific Vanquish Neo UHPLC 进行色谱分离。分离后的样品用最新一代的高分辨质谱 Orbitrap Astral 进行 DIA 质谱分析。检测模式：正离子，一级质谱扫描范围：350-1650 m/z，质谱分辨率：60,000 (@m/z 1100)，AGC target: 3e6，Maximum IT: 50 ms。MS2 采用 DIA 数据采集模式，设置 30 个 DIA 采集窗口，质谱分辨率：30,000 (@m/z 1100)，AGC target: 3e6，Maximum IT: auto，MS2 Activation Type: HCD，Normalized collision energy: 30，Spectral data type: profile。

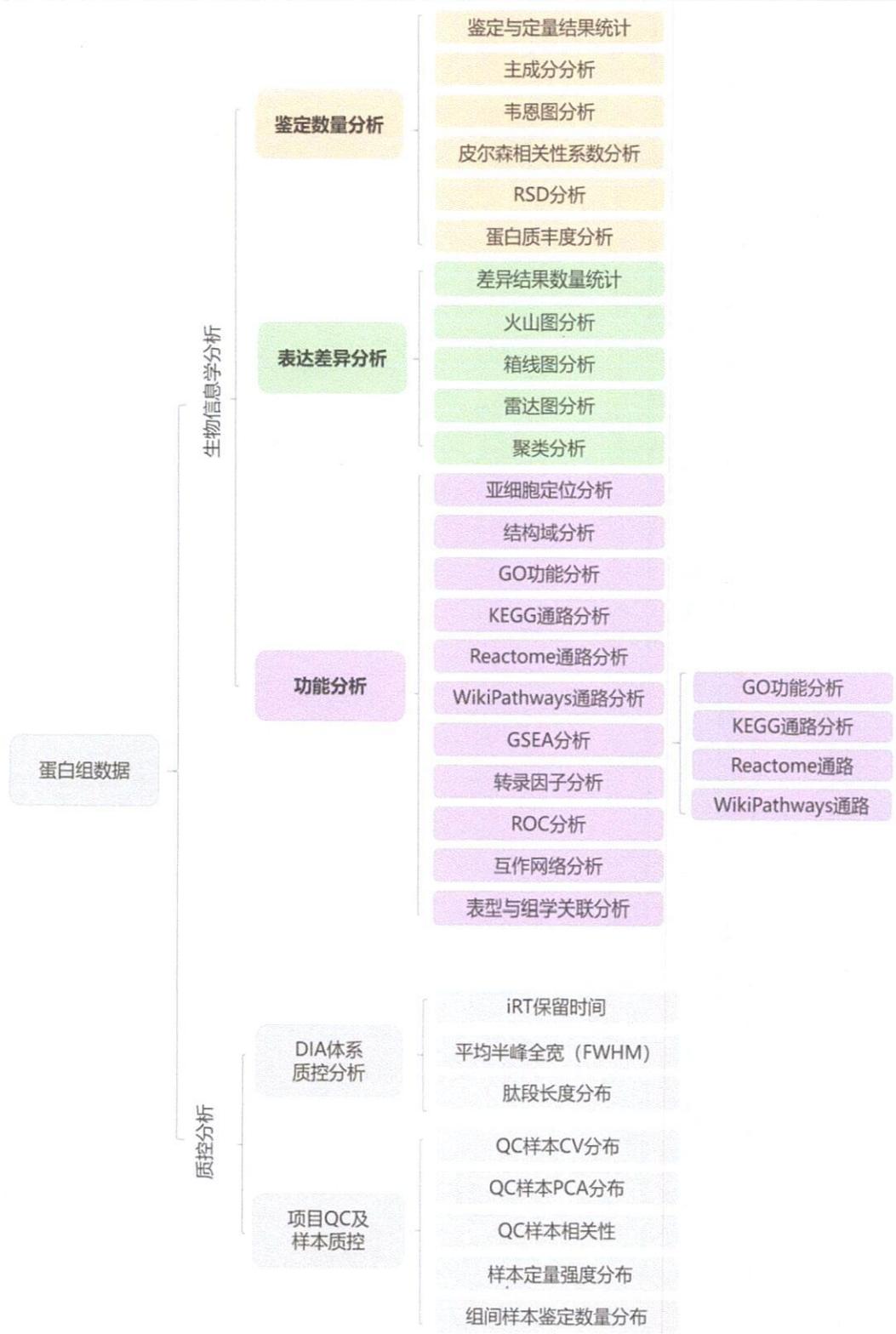
### 2.2.5 LC-MS/MS 数据分析及交付

DDA 数据直接导入商业版的 Spectronaut 软件 (Spectronaut Pulsar 18) 构建 Spectral Library。数据库采用 human\_uniprot 下载数据库。检索参数设置如下：酶为 trypsin，max miss cleavage site 为 1，固定修饰为 Carbamidomethyl(C)，动态修饰设定为 Oxidation(M) 和 Acetyl(Protein N-term)，数据库检索鉴定到的蛋白必须通过设定的过滤参数 peptide FDR <0.01, protein FDR <0.01。

DIA 数据采用 Spectronaut 软件 (Spectronaut Pulsar 18) 进行搜库和处理，软件参数设置如下：retention time prediction type 设置为 dynamic iRT, interference on MS2 level correction 为 enabled, cross run normalization 为 enabled, 所有结果必须通过设定的过滤参数 Q Value cutoff 为 0.01 (相当于 FDR<1%)。

#### **定性分析：血浆符合定性要求的鉴定数量不低于 6000**

实验下机数据搜库后进行生物信息学分析，分析内容主要包括鉴定分析、表达差异分析、功能分析等，见下图。



### 3.生物信息学分析 Bioinformatics analysis

#### 3.1 鉴定数量分析

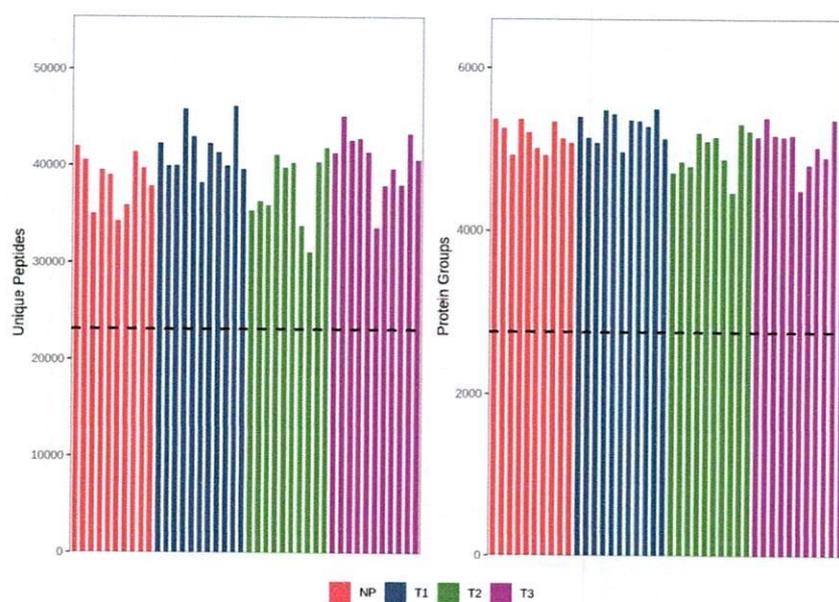
##### 3.1.1 鉴定与定量结果统计

本项目每个样本鉴定的肽段数、鉴定的蛋白数结果统计，如下表与下图。（此报告模板中仅列举部分样本鉴定数目作为示例）

DIA 鉴定与定量结果统计表

Sample	Peptides	Proteins
.....	.....	.....
.....	.....	.....
.....	.....	.....

为整体观测不同组别样本鉴定到的蛋白及肽段数目，将每个样本的鉴定结果以柱状图展示如下



DIA 鉴定结果统计柱状图

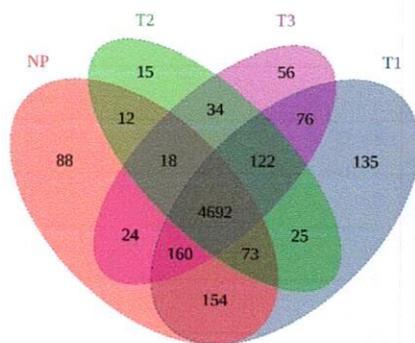
说明：Peptides：鉴定到的肽段总数目；Protein groups：鉴定到的蛋白质总数。不同颜色代表不同组别。图中虚线代表样本的蛋白/肽段鉴定数量最高的一半。

输出文件：

1) 3-1-1 鉴定与定量结果统计

### 3.1.2 组间样本鉴定重复性

为考察不同组别之间鉴定数量的重叠情况，以 Venn 图的形式将各组鉴定到的蛋白进行展示，结果如下图所示：



组间样本 Venn 图（超过五组出具花瓣图）

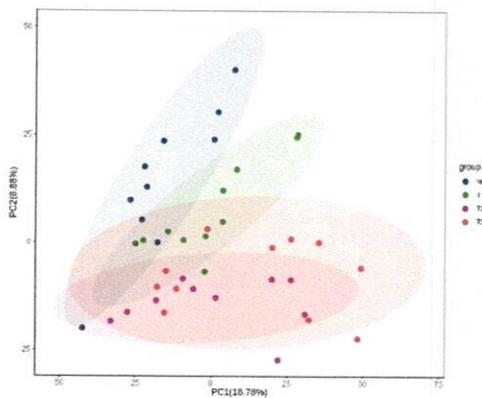
说明：每个颜色代表一个组别。

输出文件：

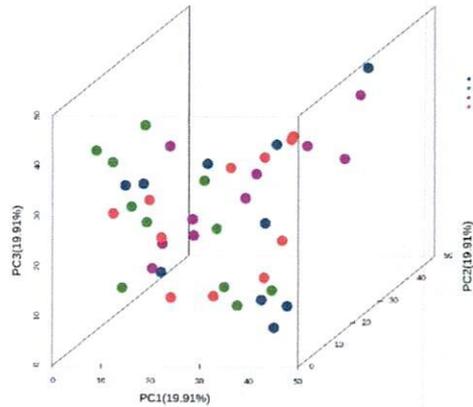
1) 3-1-2 韦恩图分析

### 3.1.3 PCA 主成分分析

主成分分析（Principal Component Analysis, PCA）是一种非监督的数据分析方法。在主成分分析中，样本的蛋白表达轮廓越相似，则聚集程度越高。样本差异越大，则距离越远，因此能从总体上反映样本组间和组内的变异性。本项目对所有样本进行 2D 和 3D PCA 分析，结果如下图所示：



所有样本 2D PCA 分布图



所有样本 3D PCA 分布图

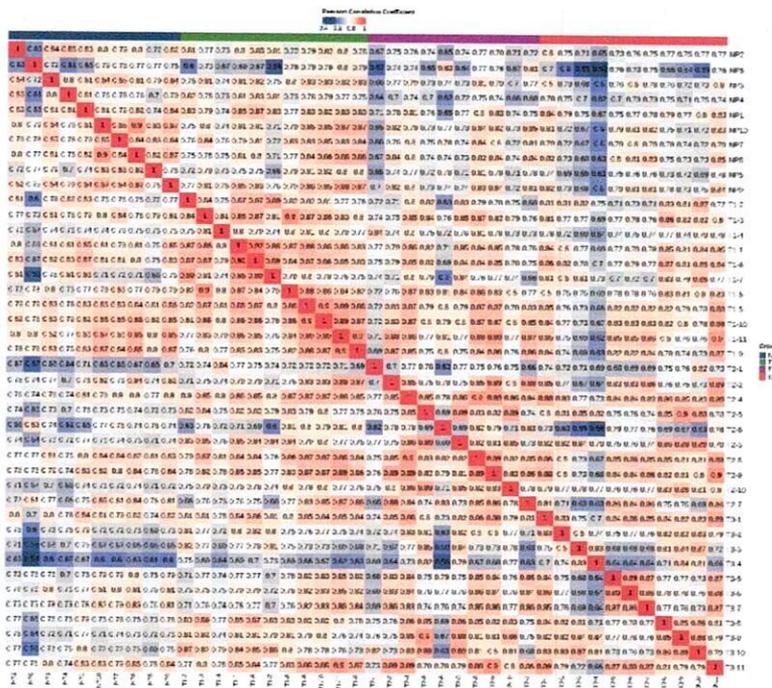
说明：图中 PC1 代表主成分 1，PC2 代表主成分 2，PC3 代表主成分 3，每个点代表一个样本，不同颜色分别代表不同的组别。

输出文件：

1) 3-1-3 PCA 分析

3.1.4 皮尔森相关性系数 (Pearson's Correlation Coefficient, PCC) 分析

所有样本两两之间计算皮尔森相关系数而绘制的热图。此系数是度量两组数据线性相关程度的值：当皮尔森系数越接近-1 为负相关，越接近 1 为正相关，越接近 0 为不相关。结果如图所示。



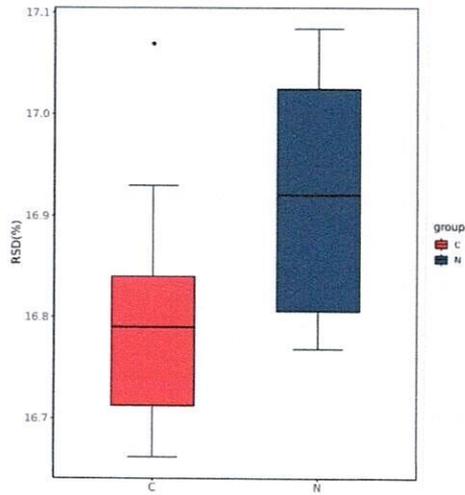
### 样本组间鉴定到的蛋白质 PCC 分析图

输出文件:

1) 3-1-4 PCC 分析

### 3.1.5 RSD 分析

样本间蛋白定量值的相对标准差 (RSD) 越小, 表明蛋白质组学的定量重复性越好。



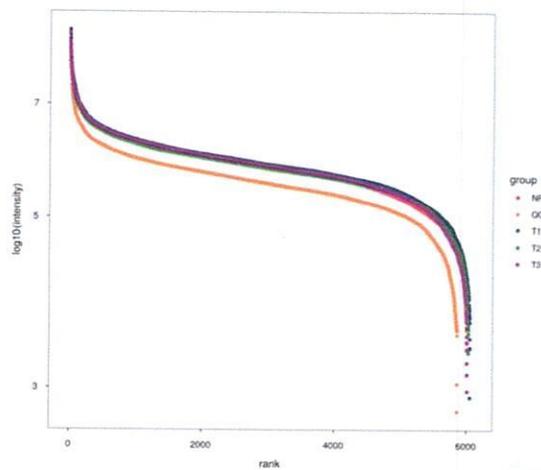
样本组间鉴定到的蛋白质 RSD 分析图

输出文件:

1) 3-1-5 BoxPlot 分析

### 3.1.6 蛋白质丰度分析

对所有组的样本鉴定到的蛋白质丰度做散点图分析, 如下所示。



蛋白质丰度分布散点图

说明: 横坐标为蛋白表达量的排名, 纵坐标为蛋白的强度值 (log10 转化)

输出文件:

### 3.2 表达差异分析

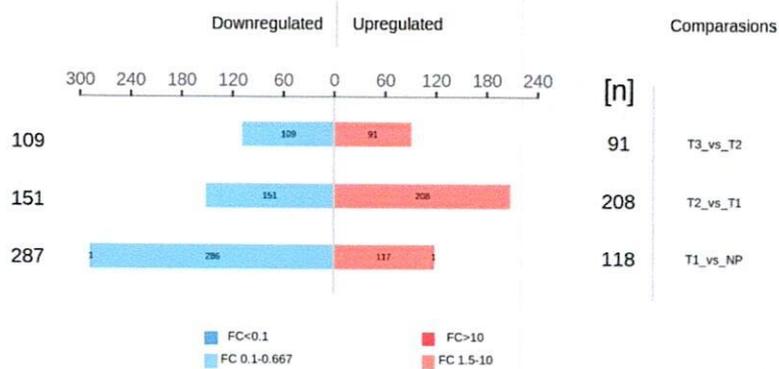
#### 3.2.1 差异结果数量统计

为了分析不同组间具有表达差异的蛋白质，对实验数据进一步进行差异筛选。

在显著性差异蛋白质筛选中，以表达倍数(Fold Change, FC) > 1.5 倍（上调大于 1.5 倍或下调小于 0.67 倍）且 P value < 0.05（T-test 或其他）为标准，得到比较组间的上调、下调蛋白质数目，如下表中 Significantly changing in abundance 列。同时，将结果以柱状图形式呈现，其中上、下调 > 10 倍的蛋白数目以更深颜色标注，如下图。

蛋白质定量差异结果统计表

Comparisons	Significantly changing in abundance			Consistent presence/absence expression profile	
	Upregulated	Downregulated	All	Upregulated	Downregulated
T3_vs_T2	109	91	1100		
T2_vs_T1	151	208	359		
T1_vs_NP	287	117	404		



蛋白质定量差异结果柱状图

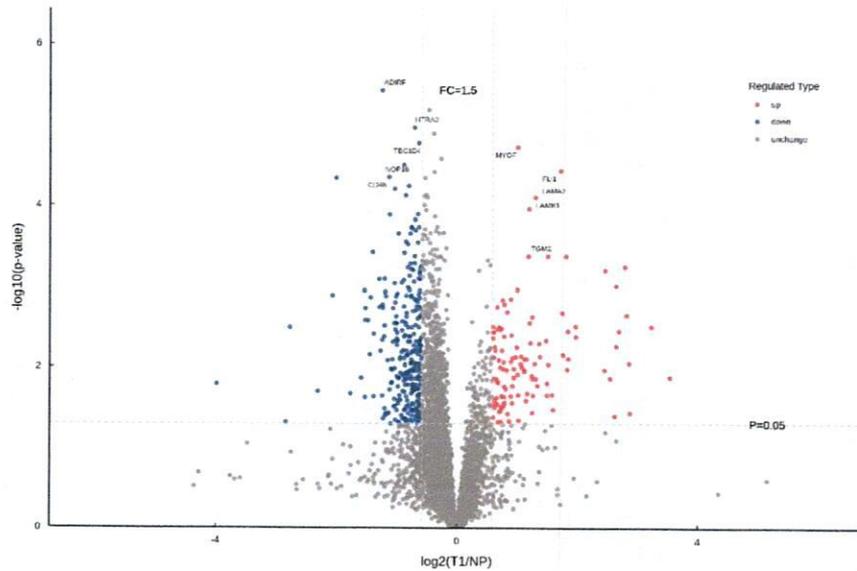
说明：Comparisons：差异比较组；Significantly changing in abundance：符合筛选倍数和 p value 的差异表达蛋白质；Consistent presence/absence expression profile：一组样品中半数及半数以上不为空值，另一组所有数据均为空值的差异蛋白质。Upregulated：上调差异表达蛋白质；Downregulated：下调表达蛋白质；All：所有差异表达蛋白质。

输出文件：

1) 3-2-1 差异结果数量统计

### 3.2.2 火山图

为了展示比较组间蛋白质的显著性差异，将比较组中蛋白质以表达差异倍数（Fold change）和 P value（T-test）两个因素为标准绘制火山图，其中显著下调的蛋白质以蓝色标注（ $FC < 0.67$  且  $p < 0.05$ ），显著上调的蛋白质以红色标注（ $FC > 1.5$  且  $p < 0.05$ ），无差异的蛋白质为灰色，并对上下调蛋白差异最显著的 top5 进行标注，结果如下图所示。



groupvs 组火山图

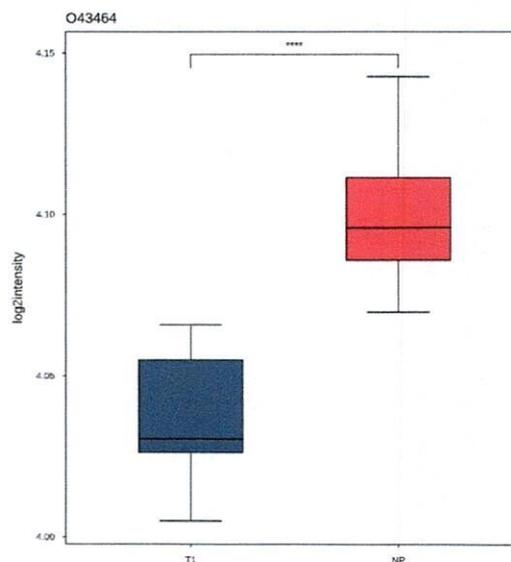
说明：横坐标为差异倍数（以 2 为底的对数变换），纵坐标为差异的显著性 P-value（以 10 为底的对数变换）。图中红色点为上调的显著性差异表达蛋白质，蓝色点为下调的显著性差异表达蛋白质，灰点为无差异变化的蛋白质。标注 ID 的点为差异最显著的 top5 上、下调蛋白。

输出文件：

1) 3-2-2 火山图

### 3.2.3 差异蛋白表达箱线图

为了更直观的展示差异蛋白在不同组之间的表达差异，利用箱线图的形式对两组间差异表达的蛋白进行展示。报告中仅展示了一个比较组一个差异蛋白的箱线图，其他比较组差异蛋白的在附件中展示箱线图。



groupvs 组差异蛋白箱线图

说明：横坐标为组别，纵坐标为表达量（以 2 为底的对数变换），图中红色和蓝色分别代表该差异蛋白在不同样本中的表达量。\*表示差异显著性程度，\*\*\*表示  $p < 0.001$ ，\*\*表示  $0.001 < p < 0.01$ ，\*表示  $0.01 < p < 0.05$ 。具体结果可查看附件。

输出文件：

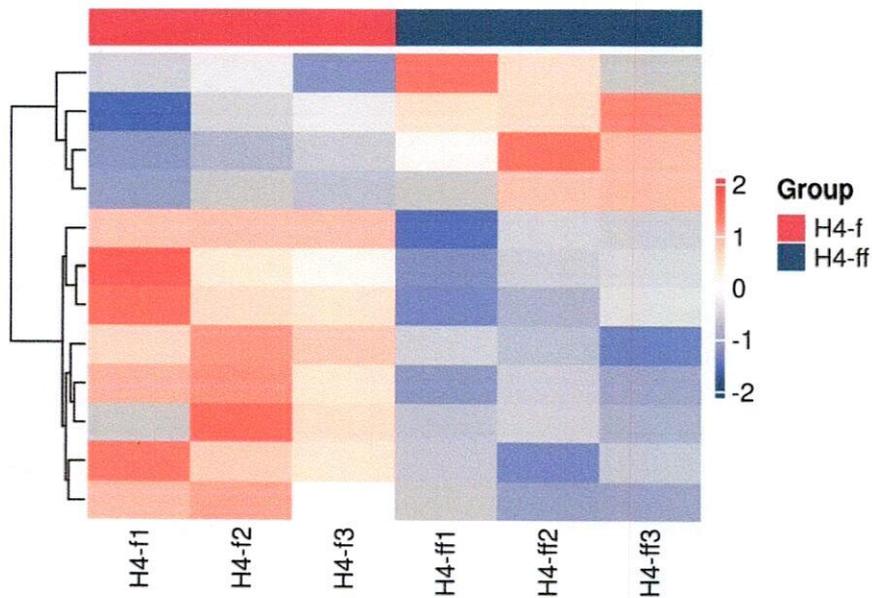
1) 3-2-3 差异蛋白箱线图

### 3.2.4 聚类分析

#### 3.2.4.1 差异蛋白表达层次聚类分析

为了分析组间、组内样本的表达模式，检验本项目分组合理性，说明差异蛋白质表达量变化是否可代表生物学处理对样本造成的显著影响，采用层次聚类算法（HierarchicalCluster）对比较组的差异表达蛋白质进行分组归类，并以热图（Heatmap）的形式展示。基于相似性基础，聚类分组结果中，一般组内的数据模式相似性较高，而组间的数据模式相似性较低，因此可以有效区分组别。

如下图，以倍数变化  $> 1.5$  倍且  $P \text{ value} < 0.05$ （T-test 或其他）的筛选标准，得到的显著差异表达蛋白质可以有效的把比较组分开，说明差异表达蛋白质筛选能够代表生物学处理对样本影响。



groupvs 组差异表达蛋白质聚类分析图

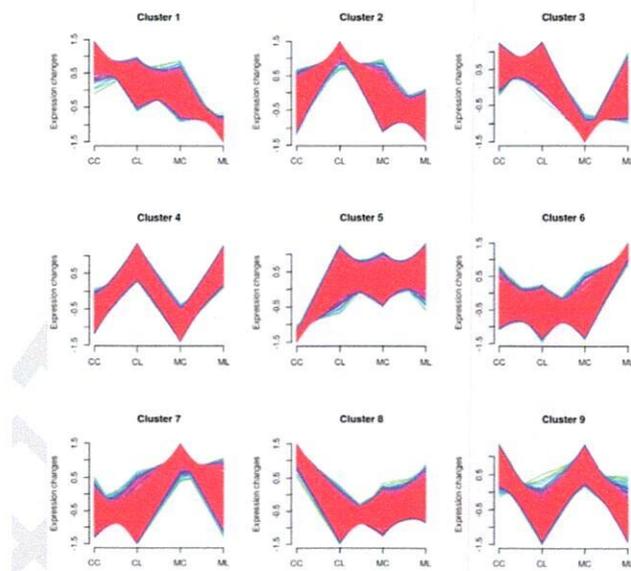
说明：层次聚类结果以树型热图表示，横轴表示样本，用不同颜色表示样本分组信息，纵轴代表差异表达的蛋白质（即纵坐标为显著性差异表达的蛋白质），显著性差异的蛋白质在不同样品中的表达量用 Z-score 方法进行标准化后以不同颜色在热图中展现，其中红色代表显著性上调的蛋白质，蓝色代表显著性下调的蛋白质，灰色部分代表无蛋白质定量信息。其他图示见附件。

输出文件：

1) 3-2-4 差异蛋白聚类分析

3.2.4.2 多组蛋白表达模式聚类

为了分析多组样本的所有蛋白整体表达模式，说明蛋白质表达量变化趋势。采用 Mfuzz 软件的 fuzzy c-means (FCM) 算法进行分析，根据所有蛋白的表达趋势分为不同的表达模块。本项目的表达模式及趋势分类如下图展示。（本分析仅适用于 3 组及以上，2 组及以下无此分析，如有具有时间梯度或者不同疾病进程阶段，需作图之前说明标注顺序）。



多组蛋白表达趋势聚类图

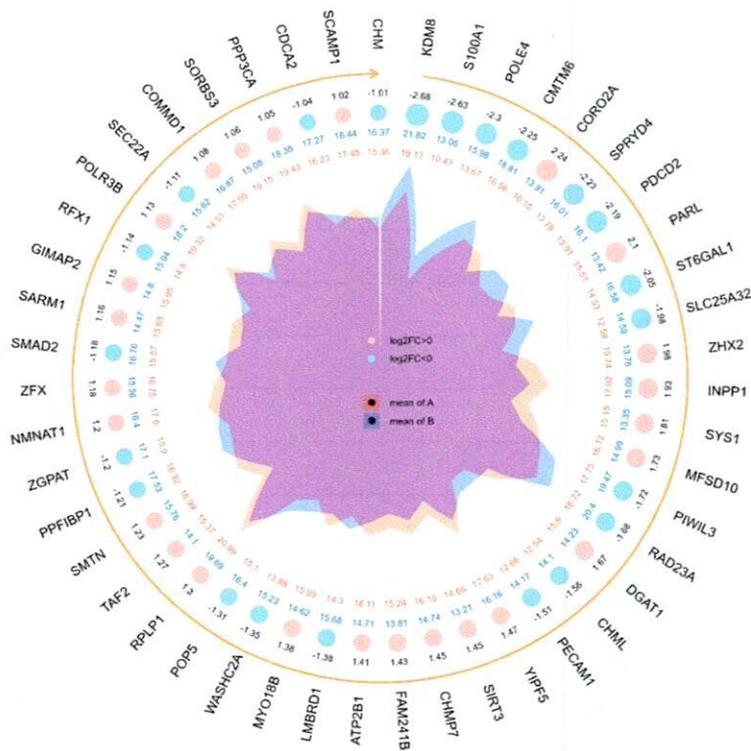
说明：横坐标代表不同组别，纵坐标表示均一化之后的表达量变化。每一个 Cluster 的线条指蛋白中表达趋势的一类蛋白。

输出文件：

- 1) 3-2-4模糊C-均值聚类分析

### 3.2.5 雷达图 (Radar Chart)

用于展示多个差异蛋白在比较组中的相对表达水平。第一圈表示多个差异蛋白（人和小鼠展示的是差异蛋白所对应的基因）；第二圈的橙色箭头当样本有重复时，表示差异蛋白对应的 P value 或者 CV 值，由小到大排序，若样本无重复，则表示差异蛋白差异倍数  $\text{Log}_2$  转换后的绝对值由大到小排序；第三圈表示  $\text{Log}_2$  转换的比较组的差异倍数，粉红色表示上调，浅蓝色表示下调，点越大表示差异倍数越大；第四圈表示两组的平均定量值。详见：



groupvs 组的差异表达雷达图

输出文件:

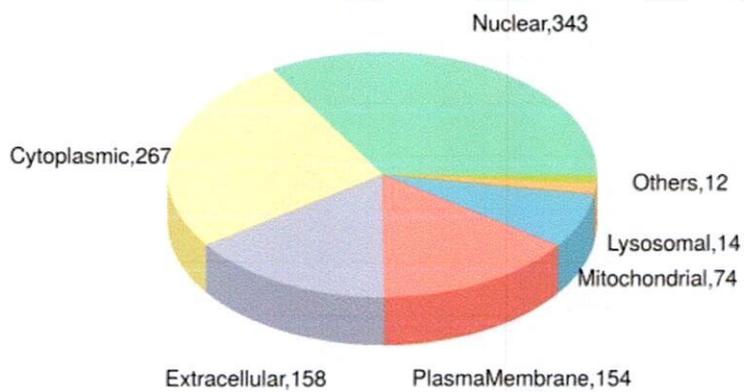
- 1) 3-2-5 差异蛋白雷达图

### 3.3 功能分析

#### 3.3.1 亚细胞定位分析

细胞器 (Organelle) 是细胞质内具有一定形态和功能的微器官 (如线粒体、内质网等), 它是蛋白发挥不同功能的重要场所。不同细胞器往往行使不同细胞功能, 故分析蛋白的亚细胞定位有助于我们进一步探究蛋白质在细胞中发挥的功能。

采用亚细胞结构预测软件 CELLO (<http://cello.life.nctu.edu.tw>)<sup>[1]</sup>对所有差异表达的蛋白质进行亚细胞定位分析, 分析结果以表格形式输出, 参见输出文件。同时, 以饼状图形式展示各细胞器中的蛋白质数目与分布比例, 如下图。

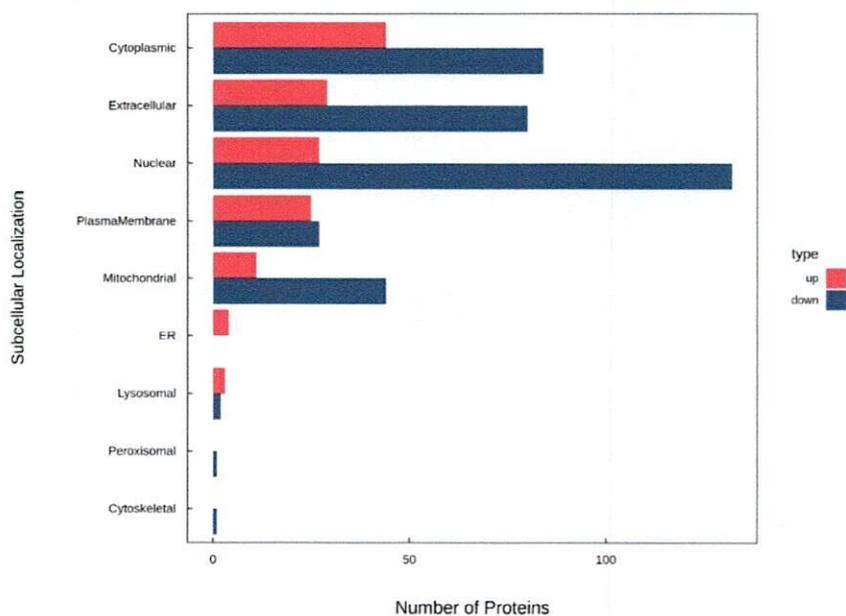


groupvs 组差异表达蛋白质亚细胞定位分布饼图

输出文件:

1) 3-3-1 亚细胞定位分析

对每个比较组的差异蛋白质进行亚细胞定位统计分析，分别统计上下调蛋白质数目，并以柱状图的形式展示；下图仅展示了一个比较组的结果，其他比较组结果见附件。



亚细胞定位结果上下调比较柱状图

说明：纵坐标代表亚细胞定位，横坐标代表该亚细胞注释到的差异蛋白质数量，红蓝色代表上下调的蛋白质。

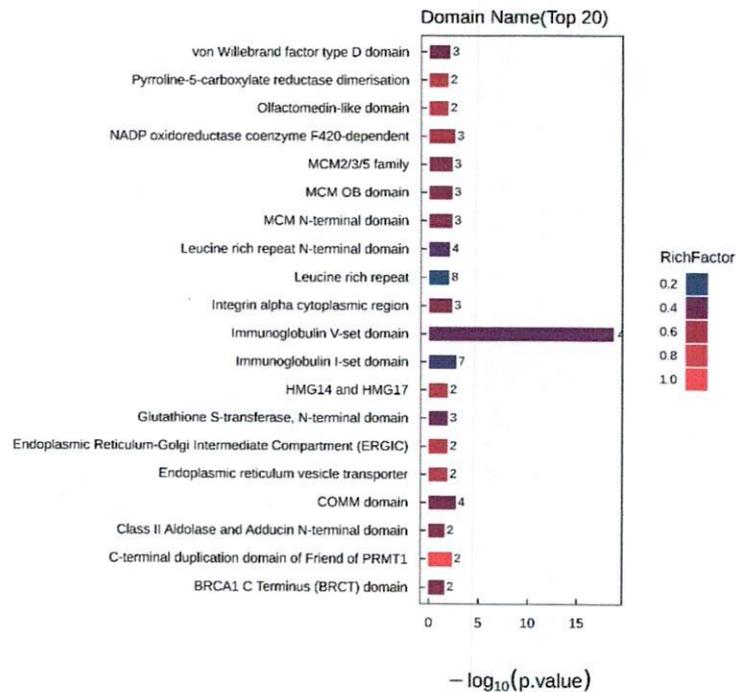
输出文件:

1) 3-3-1 亚细胞定位分析

### 3.3.2 结构域分析

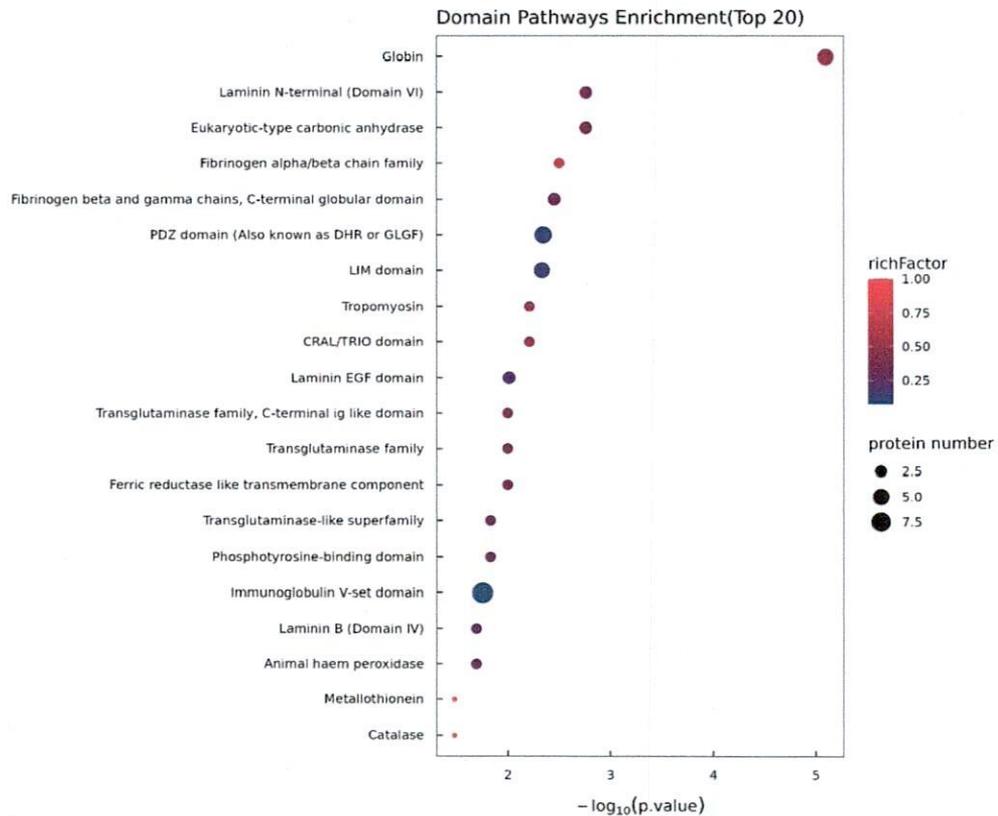
蛋白质结构域 (domain)是在较大的蛋白质分子中, 由于多肽链上相邻的超二级结构紧密联系, 形成两个或多个在空间上可以明显区别的局部区域。一般每个结构域由几十至几百个氨基酸残基组成, 各有独特的空间结构, 并承担不同的生物学功能。一般来说, 蛋白与蛋白(或其他小分子)的相互作用常以结构域为单位, 结构域内氨基酸或修饰发生改变, 可能引起蛋白关键功能的改变, 故后续氨基酸突变功能实验可以以此为参考。因此, 结构域预测对于研究蛋白关键功能区域及其发挥的潜在生物学作用具有重要意义。

采用结构域预测软件 *interproscan*<sup>[2]</sup>对差异表达蛋白质进行结构域预测, 分析结果以表格形式输出, 参见输出文件。同时, 以柱状图形式展示 Domain 中的蛋白数目(前 20), 如下图所示。



groupvs 组差异表达蛋白质结构域分析柱状图

为了揭示差异表达蛋白质的结构域富集特征, 并通过评价某个结构域条目下的蛋白质富集度的显著性水平, 找到研究者最关心的显著富集结构域及其对应差异蛋白, 采用 Fisher 精确检验 (Fisher's Exact Test) 对差异表达蛋白质进行结构域富集分析, 如下图。



groupvs 组结构域富集分析气泡图

说明：图中横坐标为某结构域分类的富集显著性，即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取  $-\log_{10}$ )，横坐标的值越大表示对应的结构域分类下富集度的显著性水平越高，颜色梯度代表富集因子的大小 (Rich Factor $\leq 1$ )，富集因子表示注释到某结构域的差异表达蛋白质数目占注释到该结构域的所有鉴定到的蛋白质数目的比例，颜色越接近红色代表 Rich Factor 值越大，气泡的大小表示每个结构域分类下差异蛋白质数目。

输出文件：

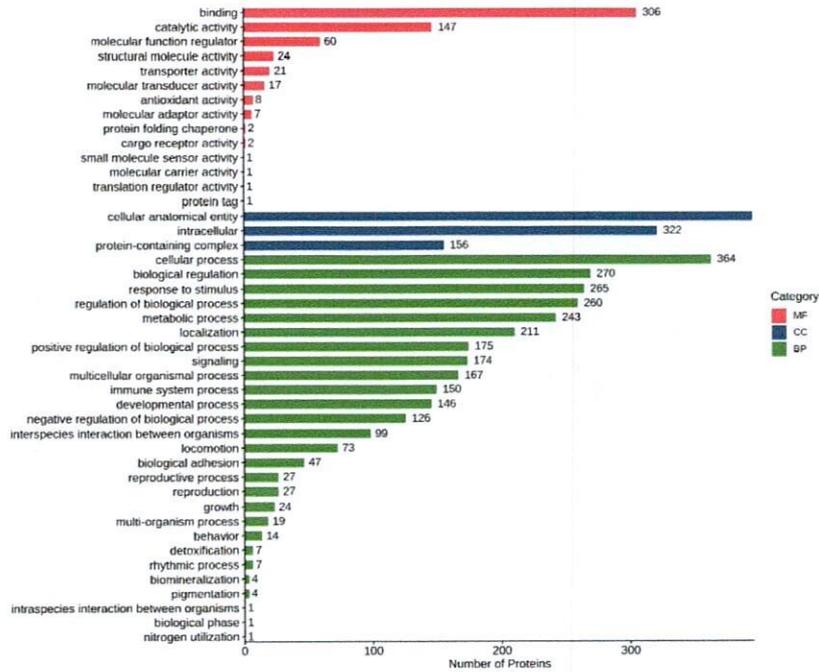
#### 1) 3-3-2 结构域分析

### 3.3.3 GO 功能分析

为了全面了解蛋白在生物体中的功能、定位及参与的生物学途径，通过基因本体 (Gene Ontology, GO) 对蛋白质进行注释。GO 是一个标准化的功能分类体系，提供了一套动态更新的标准词汇表用以描述生物体中基因和基因产物的属性。GO 功能注释主要分为 3 类：生物过程 (Biological Process, BP)，分子功能 (Molecular Function, MF) 和细胞组分 (Cellular Component, CC) [3]。

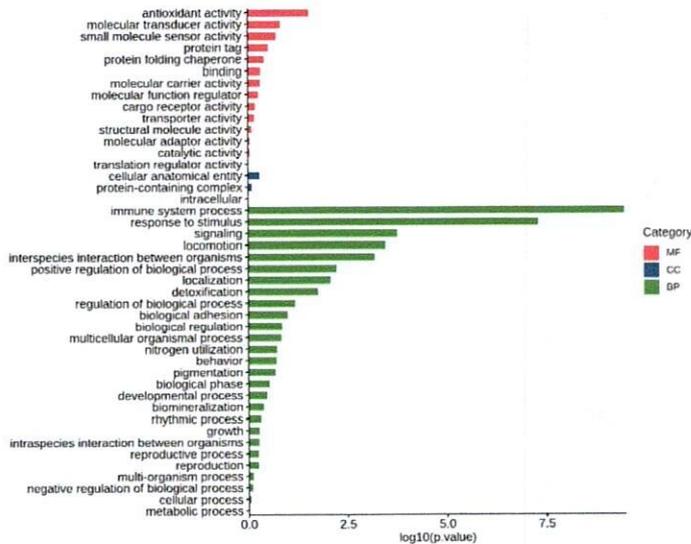
本项目采用 Blast2Go (<https://www.blast2go.com/>) [4] 软件分别对所有显著差异表达蛋白质、显著上调差异蛋白质、显著下调差异蛋白质进行 GO 功能注释，注释结果表格参见输出文件。同时，在 GO 二级功能注释层级上对显著差异蛋白数目进行统计，结果如下。

### 3.3.3.1 所有显著差异表达蛋白质 GO 功能分析



groupvs 组所有显著差异表达蛋白质的 GO 注释统计图 (level 2)

说明: 图中纵坐标表示 GO 二级功能注释信息 (GO Level2), 包含分子功能 (Molecular Function), 细胞组分 (Cellular Component) 和生物过程 (Biological Process), 依次以红色, 蓝色, 绿色予以区分; 横坐标表示每个功能分类下的显著差异表达蛋白质数目。一般情况下, 某一功能类别对应的差异表达蛋白质数目越多, 说明该功能越重要, 需要重点关注或者进行后续深入机制的探讨。



### groupvs 组所有显著差异表达蛋白质的 GO 注释统计图 (level 2)

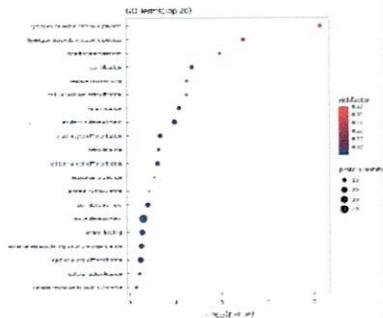
说明 图中纵坐标表示 GO 二级功能注释信息 (GO Level2), 包含分子功能 (Molecular Function), 细胞组分 (Cellular Component) 和生物过程 (Biological Process), 依次以红色, 蓝色, 绿色予以区分; 横坐标表示富集显著性, 即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取 $-\log_{10}$ ), 横坐标的值越大表示对应的 GO 功能下富集度的显著性水平越高。

输出文件:

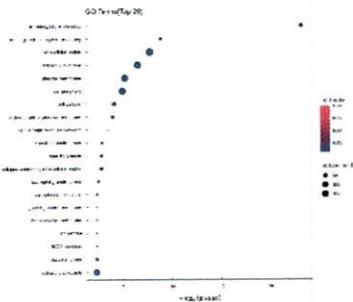
#### 1) 3-3-3 GO 功能分析

为了揭示所有差异表达蛋白质的整体功能富集特征, 并通过评价某个 GO 功能条目的蛋白质富集度的显著性水平, 找到研究者最关心的显著富集 GO 条目, 采用 Fisher 精确检验 (Fisher's Exact Test) 对差异表达蛋白质进行 GO 功能富集分析。

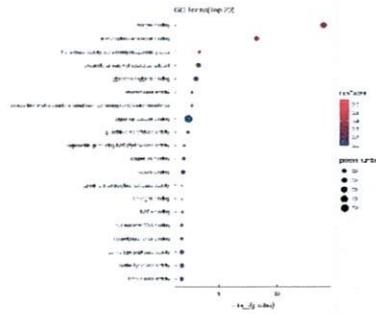
将所有差异表达蛋白质与参考物种的全部蛋白质 (或实验鉴定到的所有蛋白质) 以 GO 功能的注释结果进行对照比较, 通过 Fisher 精确检验 (Fisher's Exact Test) 得出两者差异的显著性, 从而找到所有差异表达蛋白质富集的功能类别 ( $P \text{ value} < 0.05$ )。用气泡图分别显示 GO 三大分类下的 GO 条目富集情况。



groupvs 组所有差异蛋白质的 GO 富集气泡图(BP)



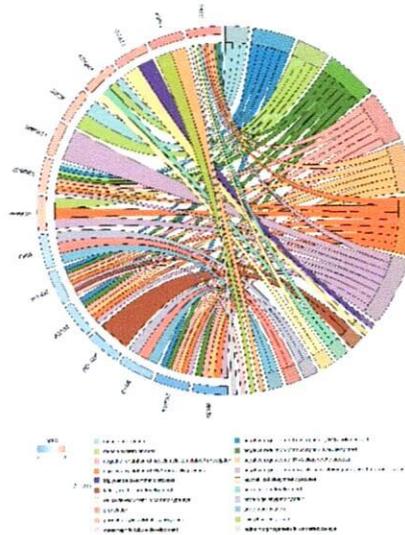
groupvs 组所有差异蛋白质的 GO 富集气泡图(CC)



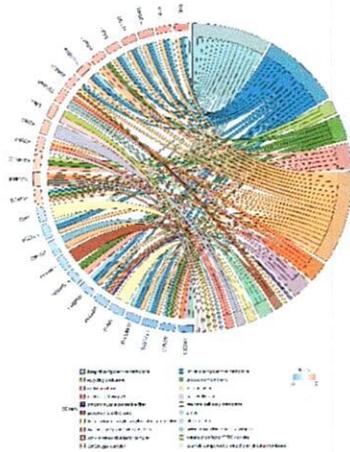
groupvs 组所有差异蛋白质的 GO 富集气泡图 (MF)

说明：图中横坐标为某 GO 功能的富集显著性，即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取  $-\log_{10}$ )，横坐标的值越大表示对应的 GO 功能下富集度的显著性水平越高，颜色梯度代表富集因子的大小 ( $\text{Rich Factor} \leq 1$ )，富集因子表示注释到某 GO 功能的差异表达蛋白质数目占注释到该 GO 功能的所有鉴定到的蛋白质数目的比例，颜色越接近红色代表 Rich Factor 值越大，气泡的大小表示每个 GO 功能分类下差异蛋白质数目。一般情况下，GO 富集结果中 P 值越小 ( $P < 0.05$ )，对应 GO 功能分类从统计学上讲富集越显著，而与 GO 功能分类相关的差异表达蛋白质数目在某种程度上反映实验设计中生物学处理对各个分类的影响程度大小，因此可以结合两方面因素，选择较为感兴趣的生物学功能以及显著性影响这些功能的差异表达蛋白质进行后续生物学实验验证或机制研究。

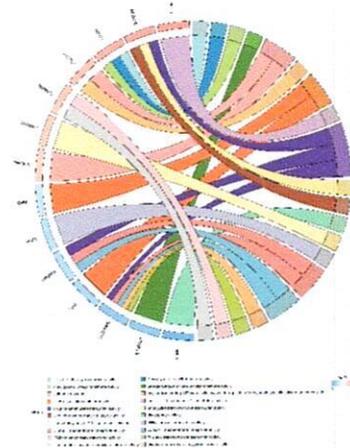
显著富集弦图，用于展示显著富集的 GO 功能与蛋白之间的关系，图的右侧表示富集到的 GO 功能，与右侧功能相连的是该功能中的差异蛋白，差异蛋白的顺序依据其  $\text{Log}_2\text{FC}$  值从大到小排列。该图能直观的展示富集 GO 功能中每个蛋白的名字、差异程度。



groupvs 组所有差异蛋白质的 GO 富集弦图 (BP)



groupvs 组所有差异蛋白质的 GO 富集弦图(CC)

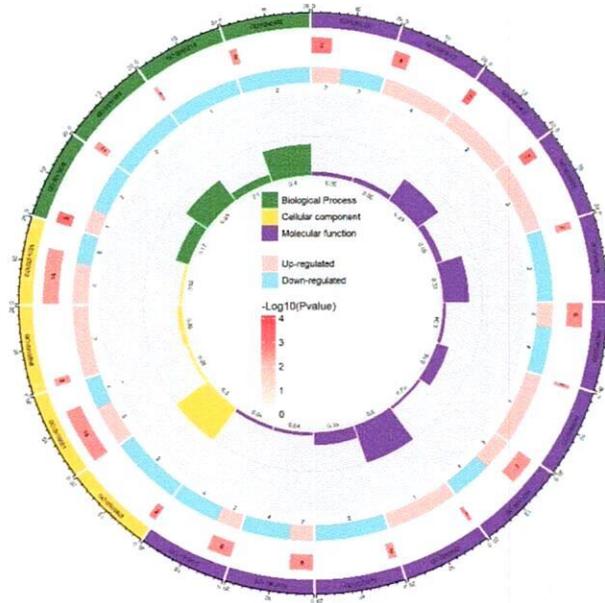


groupvs 组所有差异蛋白质的 GO 富集弦图(MF)

输出文件：

1) 3-3-3 GO 功能分析

Circos 每一圈含义（由外到内）：第一个圆圈：富集的 GO 功能一级分类，圆圈外是蛋白数量的坐标标尺。不同的颜色代表不同的类别；第二个圆圈：功能富集显著性 P value 经-Log10 转换后的值。数值越大，颜色越红；第三圈：上、下调差异蛋白数量条形图，红色代表上调差异蛋白数量，蓝色代表下调差异蛋白数量；第四个圆圈：每个功能的富集因子的大小（Rich Factor $\leq$ 1）。注：富集条目少于 4 不显示。



groupvs 组所有差异蛋白质的 GO 富集 Circos 图

输出文件:

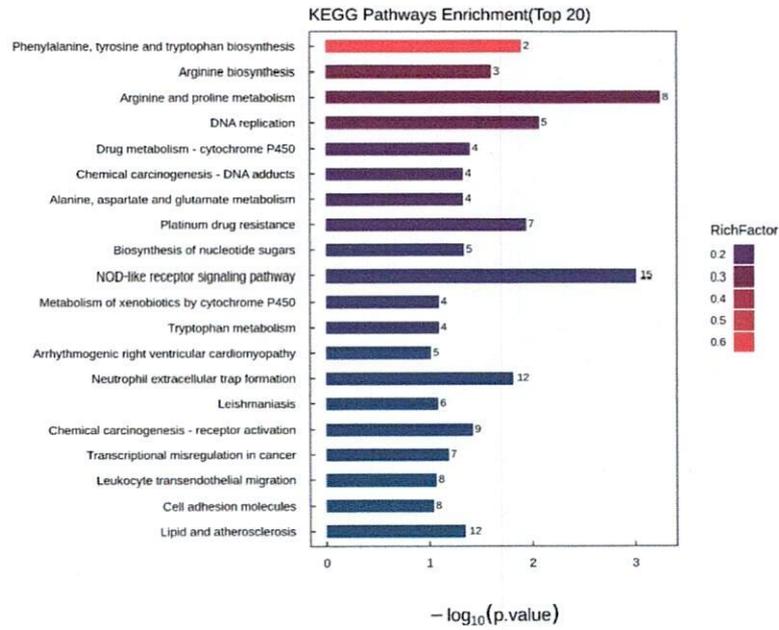
1) 3-3-3 GO 功能分析

### 3.3.4 KEGG 通路分析

为了更系统全面地解析生物学过程、疾病发生机理、药物作用机制等，往往需要从一系列蛋白质协调作用的角度阐述变化规律，如代谢通路变化。因此通过 KEGG (Kyoto Encyclopedia of Genes and Genomes) 数据库对蛋白质解析注释<sup>[5]</sup>。KEGG 是由研究人员阅读海量文献后，将众多的代谢途径以特定的图形语言整理而成的数据库，其收录了新陈代谢，遗传信息加工，环境信息加工，细胞过程，生物体系统，人类疾病以及药物开发等多个方面的通路信息，常用于通路研究。

#### 3.3.4.1 所有显著差异蛋白质 KEGG 通路注释分析

本项目将所有显著差异蛋白质进行 KEGG 通路注释，注释表格参见输出文件。结果如下图所示。更多信息请参考: <http://www.genome.jp/kegg/pathway.html>。



groupvs 组显著差异蛋白质的 KEGG 通路注释统计图 (Top20)

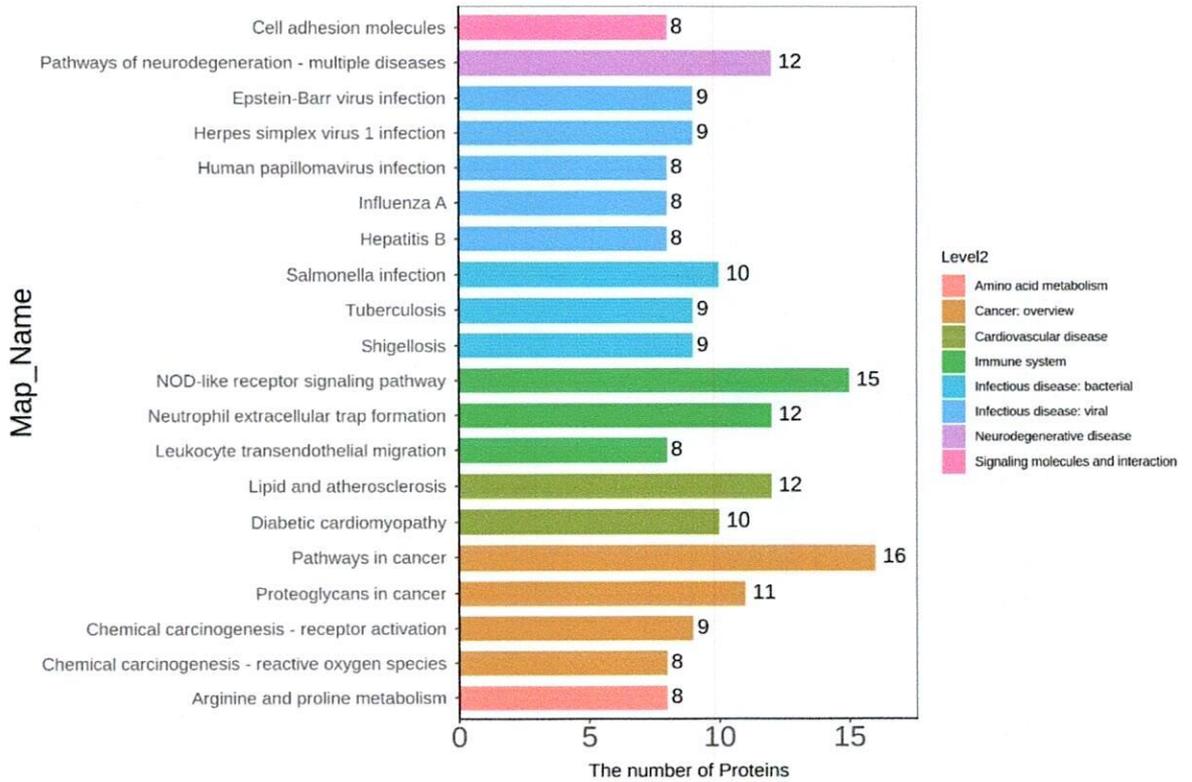
说明：图中纵坐标是含差异表达蛋白质参与的通路名称，横坐标表示富集显著性，即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取  $-\log_{10}$ )，横坐标的值越大表示对应的通路下富集度的显著性水平越高。一般情况下，参与某一通路的差异表达蛋白质数目越多，说明该通路越重要，需要重点关注或者进行后续深入机制的探讨。

输出文件：

#### 1) 3-3-4KEGG 通路注释分析

### 3.3.4.2 所有显著差异蛋白质 KEGG 通路归属分析

在生物体内，不同蛋白相互协调行使其生物学行为，基于 Pathway 的分析有助于更进一步了解其生物学功能 Pathway 富集不同层级结果。KEGG 代谢通路共分为 7 个分支：细胞过程(Cellular Processes)、环境信息处理(Environmental Information Processing)、遗传信息处理(Genetic Information Processing)、人类疾病(Human Diseases)(仅限动物)、代谢(Metabolism)、有机系统(Organismal Systems)、药物开发(Drug Development)。本分析通过图形展示众多的代谢途径以及通路归属关系，以便更加直观观测到差异表达蛋白所参与的代谢途径。差异表达蛋白质参与的通路代谢注释如下展示。



groupvs 组显著差异蛋白质的 KEGG 通路注释及归属柱状图

说明：横坐标轴代表通路蛋白注释数目，纵坐标代表 KEGG 注释名称。不同颜色代表不同 KEGG 的代谢通路 level2 层级。。

输出文件：

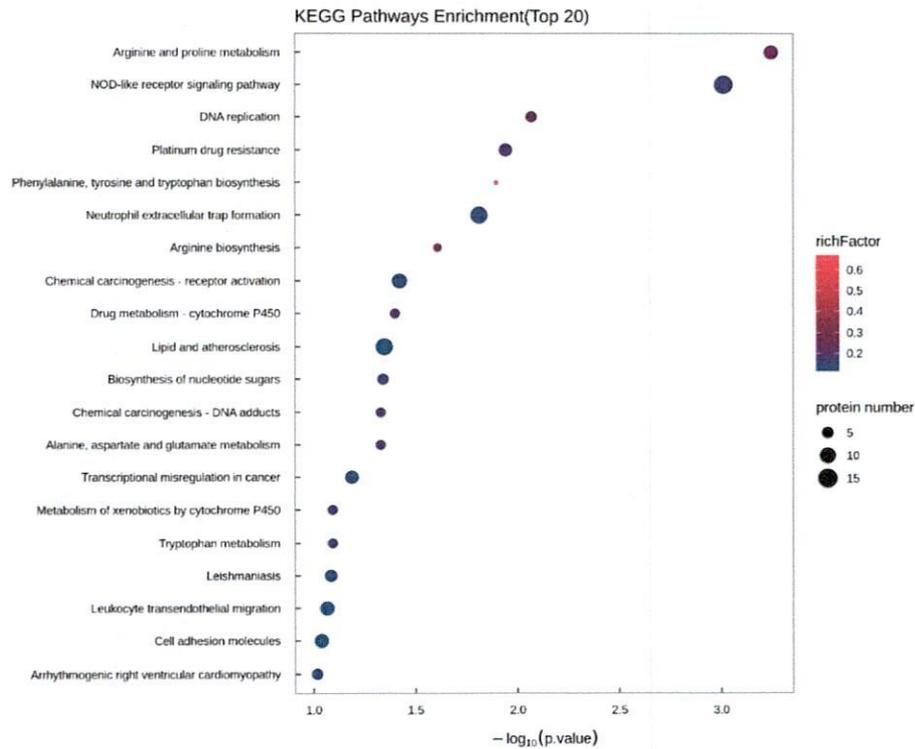
1) 3-3-4-1KEGG 通路注释分析

### 3.3.4.3 所有显著差异蛋白质 KEGG 通路富集分析

为了揭示所有差异蛋白质的整体代谢通路富集特征，并通过评价某个 KEGG 代谢通路的蛋白质富集度的显著性水平，找到研究者最关心的显著富集 KEGG 代谢通路，采用 Fisher 精确检验

(Fisher's Exact Test) 对差异表达蛋白质进行 KEGG 通路富集分析。

将所有显著差异蛋白质与参考物种的全部蛋白质（或实验鉴定到的所有蛋白质）以 KEGG 的注释结果进行对照比较，通过 Fisher 精确检验 (Fisher's Exact Test) 得出两者差异的显著性，从而找到所有差异表达蛋白质富集的通路类别 (P value < 0.05)。如下图所示，通过 Fisher 精确检验方法对 groupvs 比较组的显著差异蛋白质进行 KEGG 通路富集分析。



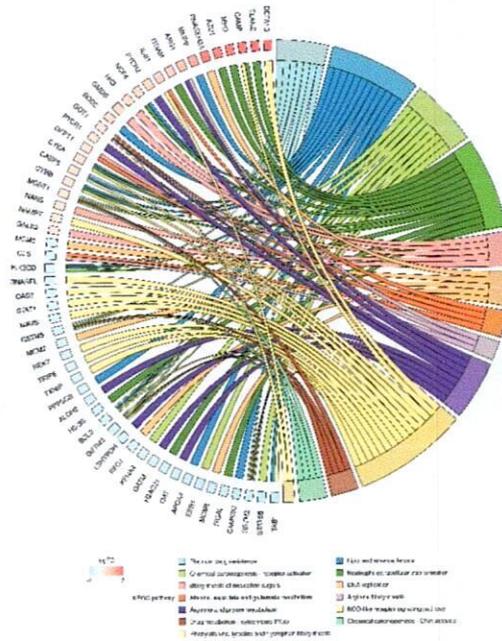
groupvs 组所有差异蛋白质的 KEGG 通路富集气泡图

说明：图中横坐标为某 KEGG 通路的富集显著性，即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取  $-\log_{10}$ )，横坐标的值越大表示对应代谢通路富集度的显著性水平越高，颜色梯度代表富集因子的大小 (Rich Factor  $\leq 1$ )，富集因子表示注释到 KEGG 通路类别的显著差异表达蛋白质数目占注释到该类别的所有鉴定到的蛋白质数目的比例，颜色越接近红色代表 Rich Factor 值越大，气泡的大小表示每个 KEGG 通路下差异蛋白质数目。因此可以选择较为感兴趣的生物学功能以及显著性影响这些功能的差异表达蛋白质进行后续生物学实验验证或机制研究。

输出文件：

### 1) 3-3-4KEGG通路富集分析

显著富集弦图，用于展示显著富集的 KEGG 通路与蛋白之间的关系，图的右侧表示富集到的 KEGG 通路，与右侧通路相连的是该通路中的差异蛋白，差异蛋白的顺序依据其 Log2FC 值从大到小排列。该图能直观的展示富集通路中每个蛋白的名字、差异程度。

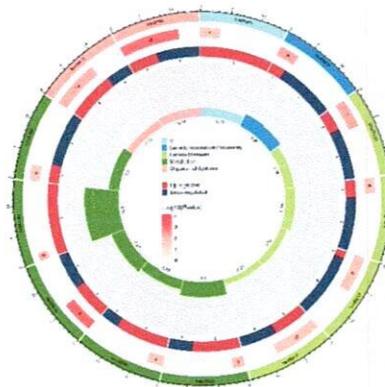


KEGG 通路的富集统计弦图

输出文件：

1) 3-3-4KEGG 通路富集分析

对差异蛋白 KEGG 通路注释结果进行富集分析，以 Circos 图形式来进行结果展示。Circos 每一圈含义（由外到内）：第一个圆圈：富集的 KEGG 通路，圆圈外是蛋白数量的坐标标尺；第二个圆圈：蛋白 KEGG 通路富集显著性 P value 经-Log10 转换后的值。数值越大，颜色越红；第三圈：上、下调差异蛋白数量条形图，红色代表上调差异蛋白数量，蓝色代表下调差异蛋白数量；第四个圆圈：每个 KEGG 通路的富集因子的大小（Rich Factor $\leq 1$ ）。注：富集条目少于 4 不显示。



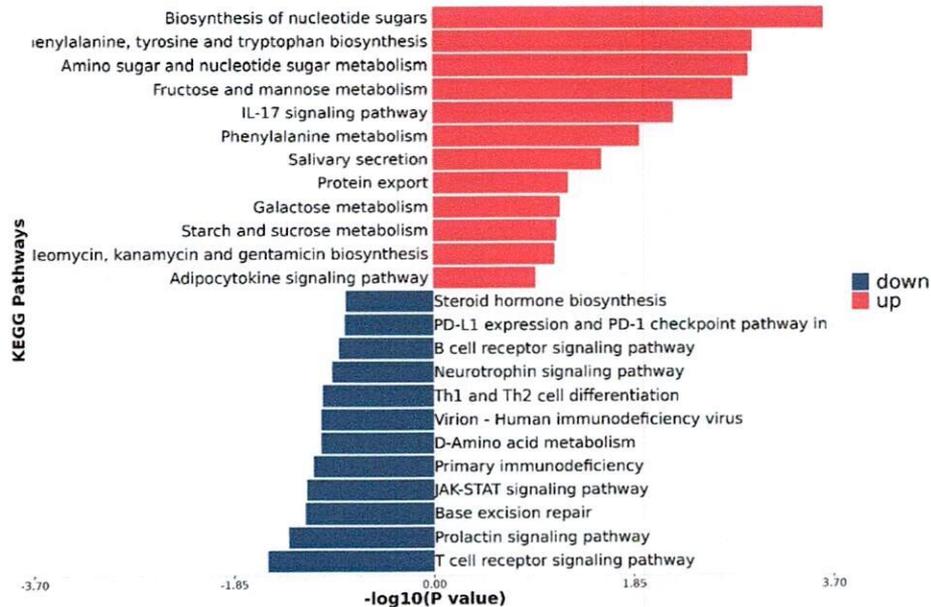
KEGG 通路的富集统计 Circos 图 (top 20)

输出文件：

- 1) [3-3-4KEGG通路富集分析](#)

### 3.3.4.4 显著上、下调差异蛋白质 KEGG 通路富集分析

为了更好地考察差异蛋白的通路富集的显著性，分别对上、下调差异表达蛋白质进行 KEGG 通路富集分析，以蝴蝶图形式展示，结果展示如下：



groupvvs 组显著上、下调差异表达蛋白质的通路富集蝴蝶图

说明：横坐标为 Fisher 精确检验的 p value 值（取以 10 为底的对数），纵坐标表示通路名称。上调和下调蛋白参与的通路用红色（右）和蓝色（左）条表示。

每个差异蛋白注释到的通路图信息及在通路中的位置，在附件中展示。

输出文件：

- 1) [3-3-4KEGG 通路富集分析](#)

每个差异蛋白注释到的通路图信息及在通路中的位置，在附件中展示。

输出文件：

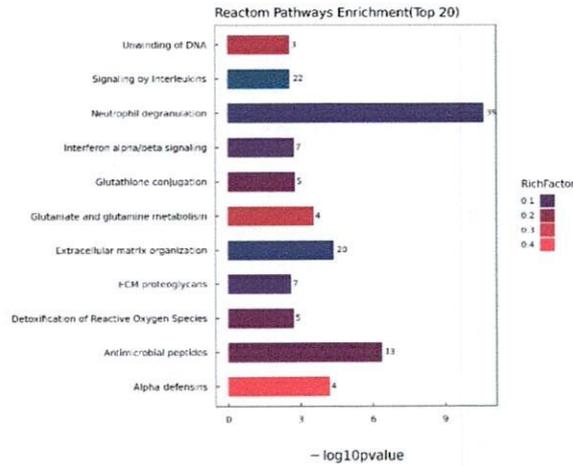
- 1) [3-3-4KEGG 通路富集分析](#)

### 3.3.5 Reactome 通路富集分析

Reactome 是一个免费提供的开源关系数据库，其中包含信号和代谢分子及其组织成生物途径和过程的关系。Reactome 数据模型的核心单元是反应。参与反应的实体（核酸、蛋白质、复合物、疫

苗、抗癌治疗剂和小分子) 形成生物相互作用网络, 并被分组为通路。Reactome 中的生物学途径的例子包括经典的中间代谢、信号传导、转录调控、细胞凋亡和疾病。限人、大鼠、小鼠物种。

显著富集条形图, 横轴表示 $-\log_{10}$  转换的富集显著性 P value, 纵轴为对应 Reactome 通路描述信息。条形图长短表示富集显著性, 越长代表富集显著性越强。

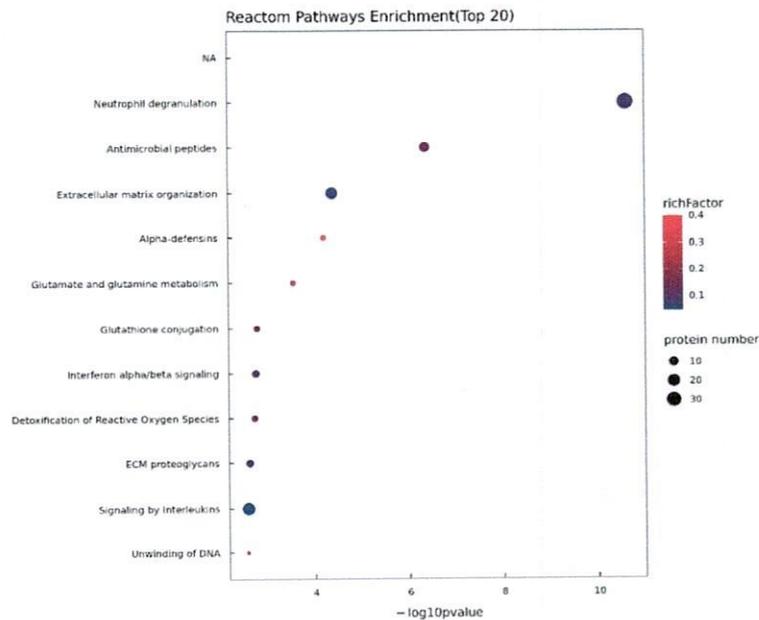


groups 组所有差异蛋白质的富集条形图

输出文件:

### 1) 3-3-5 Reactome 通路富集分析

给出了最显著富集的前 20 个功能的结果, 图中纵轴为 Reactome 通路描述信息, 横坐标为某通路的富集显著性, 即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取 $-\log_{10}$ ), 横坐标的值越大表示对应代谢通路富集度的显著性水平越高, 颜色梯度代表富集因子的大小 (Rich Factor  $\leq 1$ ), 富集因子表示注释到该通路类别的显著差异表达蛋白质数目占注释到该类别的所有鉴定到的蛋白质数目的比例, 颜色越接近红色代表 Rich Factor 值越大, 气泡的大小表示该通路下差异蛋白质数目。

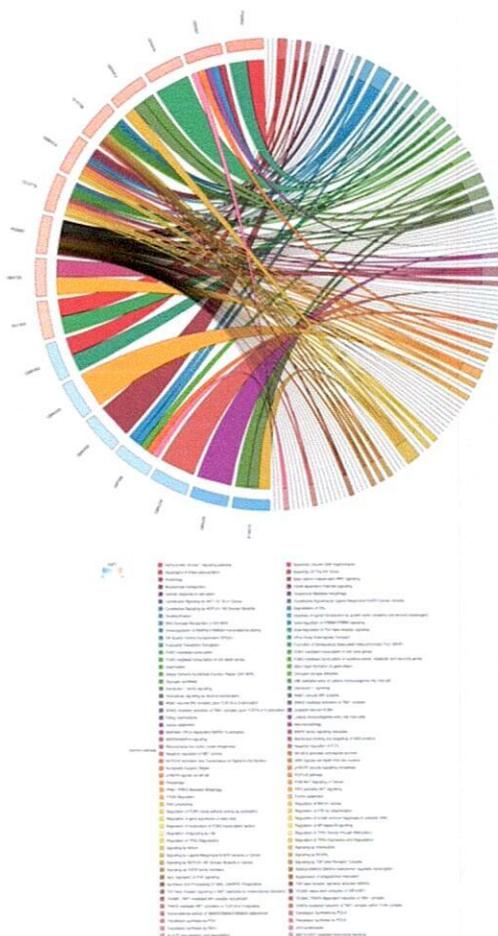


## groupvs 组所有差异蛋白质的富集气泡图

输出文件:

### 1) 3-3-5 Reactome 通路富集分析

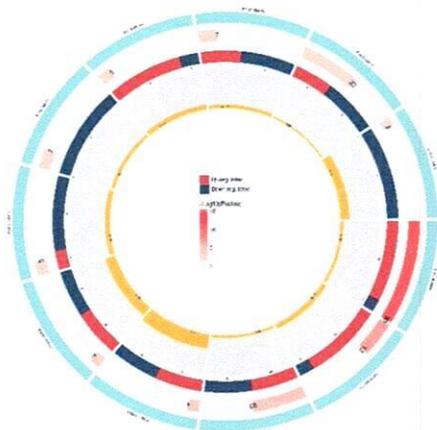
显著富集弦图, 用于展示显著富集的 Reactome 通路和蛋白之间的关系, 图的右侧表示富集到的 Reactome 通路, 与右侧 Reactome 通路相连的是该通路中的差异蛋白, 差异蛋白的顺序依据其 Log2FC 值从大到小排列。该图能直观的展示富集通路中每个蛋白的名字、差异程度。



输出文件:

### 1) 3-3-5 Reactome 通路富集分析

Circos 每一圈含义 (由外到内): 第一个圆圈: 富集的 Reactome 通路, 圆圈外是蛋白数量的坐标标尺; 第二个圆圈: Reactome 通路富集显著性 P value 经  $-\log_{10}$  转换后的值。数值越大, 颜色越红; 第三圈: 上、下调差异蛋白数量条形图, 红色代表上调差异蛋白数量, 蓝色代表下调差异蛋白数量; 第四个圆圈: 每个功能的富集因子的大小 ( $\text{Rich Factor} \leq 1$ )。



groupvs 组所有差异蛋白质富集 Circos 图

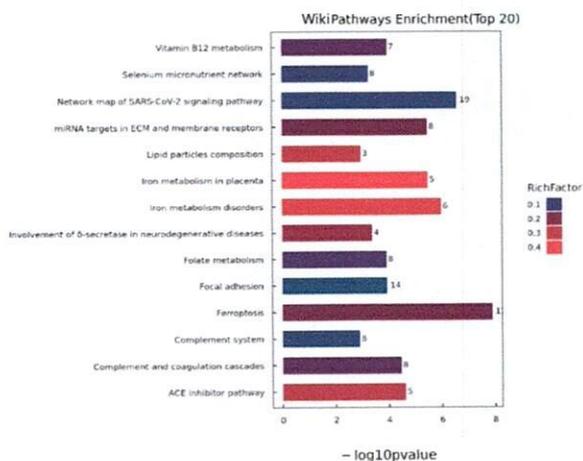
输出文件:

1) 3-3-5 Reactome 通路富集分析

**3.3.6 WikiPathways 通路富集分析**

WikiPathways 的建立是为了促进生物学界对通路信息的贡献和维护。WikiPathways 是一个开放的致力于管理生物通路数据库。因此, WikiPathways 为生物途径数据库提供了一种新模型, 可增强和补充现有通路数据库信息, 例如 KEGG、Reactome 和 Pathway Commons 等。限人、大鼠、小鼠物种。

显著富集条形图, 横轴表示 $-\log_{10}$  转换的富集显著性 P value; 纵轴为对应 WikiPathways 通路描述信息。条形图长短表示富集显著性, 越长代表富集显著性越强。

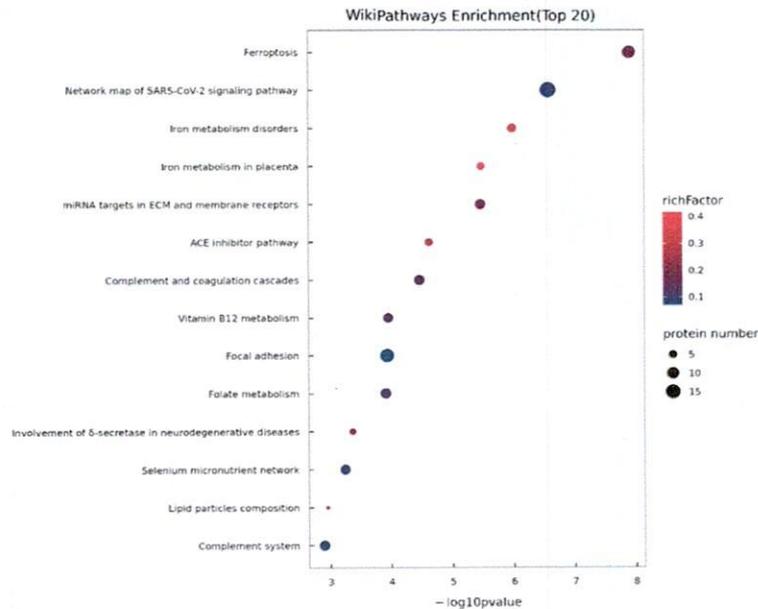


groupvs 组所有差异蛋白质的富集条形图

输出文件:

1) 3-3-6 WikiPathways 通路富集分析

显著富集气泡图中给出了最显著富集的前 20 个功能的结果，图中纵轴为 WikiPathways 通路描述信息，横坐标为某通路的富集显著性，即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取 $-\log_{10}$ )，横坐标的值越大表示对应代谢通路富集度的显著性水平越高，颜色梯度代表富集因子的大小 (Rich Factor $\leq 1$ )，富集因子表示注释到该通路类别的显著差异表达蛋白质数目占注释到该类别的所有鉴定到的蛋白质数目的比例，颜色越接近红色代表 Rich Factor 值越大，气泡的大小表示该通路下差异蛋白质数目。

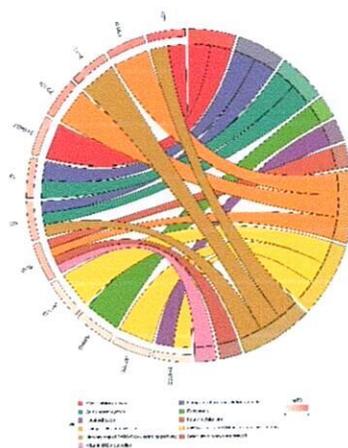


groupvs 组所有差异蛋白质的富集气泡图

输出文件:

### 1) 3-3-6 WikiPathways 通路富集分析

显著富集弦图，用于展示显著富集的 WikiPathways 通路蛋白之间的关系，图的右侧表示富集到的 WikiPathways 通路，与右侧 WikiPathways 通路相连的是该通路中的差异蛋白，差异蛋白的顺序依据其 Log2FC 值从大到小排列。该图能直观的展示富集通路中每个蛋白的名字、差异程度。

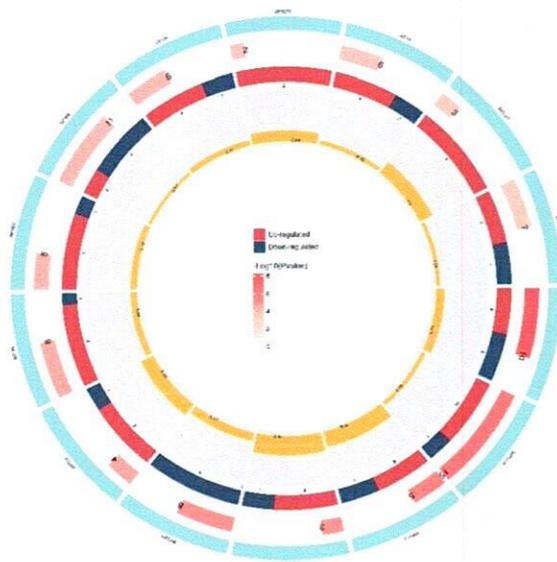


## groupvs 组所有差异蛋白质富集弦图

输出文件:

### 1) 3-3-6 WikiPathways 通路富集分析

Circos 每一圈含义 (由外到内): 第一个圆圈: 富集的 WikiPathways 通路; 第二个圆圈: WikiPathways 通路富集显著性 P value 经-Log10 转换后的值。数值越大, 颜色越红; 第三圈: 上、下调差异蛋白数量条形图, 红色代表上调差异蛋白数量, 蓝色代表下调差异蛋白数量; 第四个圆圈: 每个功能的富集因子的大小 (Rich Factor $\leq 1$ )。注: 富集条目少于 4 不显示。



groupvs 组所有差异蛋白质富集 Circos 图

输出文件:

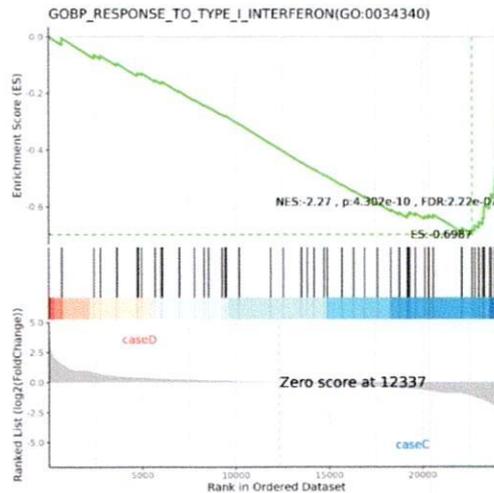
### 1) 3-3-6 WikiPathway 通路富集分析

#### 3.3.7 GSEA 分析

传统的蛋白功能富集方法是基于超几何检验而针对差异表达蛋白进行富集的, 但当单个蛋白表达变化不大时, 基于传统富集分析得到结果可能会很少, 甚至没有结果。GSEA 分析 (Gene Set Enrichment Analysis) 能够有效弥补传统富集分析对信息挖掘不足等问题, 更能全面地对某一功能单位 (通路、GO term 或其他) 的调节作用进行解释。它可以将那些在传统富集分析信息中容易遗漏掉的差异表达不显著却有着重要生物学意义的基因包含在内, 也可以解决传统富集分析中因为得到的差异基因较少, 而无法开展功能富集分析或者无法富集到感兴趣的通路的问题。其基本思想是使用预定义的蛋白集, 将蛋白按照在两类样本中的差异表达程度排序, 然后检验预先设定的蛋白集是否在这个排序表的顶端或者末端富集。此外, 通过 GSEA 分析可以判断某条通路中基因的总变化趋势, 以及该通路到底是激活还是抑制状态。GSEA 富集分析主要包括三个步骤: 计算富集得分 (Enrichment Score); 估计富集得分的显著性水平; 矫正多重假设验证。我们分别对物种的 GO、KEGG、Reactome、WikiPathways 数据集进行 GSEA 分析, 显著富集的蛋白集呈图展示见附件。限人、大鼠、小鼠、果蝇和酵母物种。

### 3.3.7.1 GSEA GO 功能分析

通过 GSEA 的方式将可定量蛋白进行 GO 功能条目富集分析。



Groupvs 组 GSEA GO 富集分析

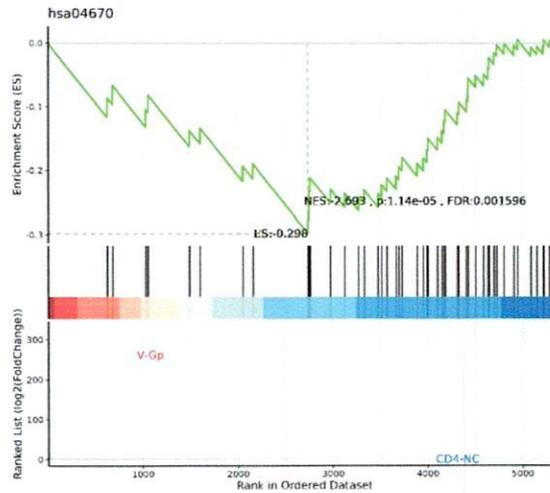
说明：横轴为比较组中的蛋白（蛋白集）根据其表达量变化值（Log2FC）由大到小的排序。GSEA 富集图自上而下分为三部分：①上部分显示的是当分析沿着蛋白集按排序计算时，ES（Enrichment Score）值在计算到每个蛋白位置时的展示（即分析过程中动态的 ES 值），最高峰处的 ES 得分即为该通路的 ES 值。②中间部分俗称条形码图，用线条标记了该通路中涉及到的蛋白出现在蛋白集排序列表中的位置。红蓝相间的热图是表达丰度排列（红色越深的表示该位置的基因 log2FC 越大，蓝色越深表示 log2FC 越小）。③最下面部分为排序后比较组中蛋白 Log2FC 的排序，以灰色面积图展示。图的右上侧的注释为 GO 通路富集 pvalue 和 FDR 值。

输出文件：

1) 3-3-7 GSEA GO 分析

### 3.3.7.2 GSEA KEGG 富集分析

通过 GSEA 的方式将可定量蛋白进行 KEGG 通路富集分析。



### Groupvs 组 GSEA KEGG 通路富集分析

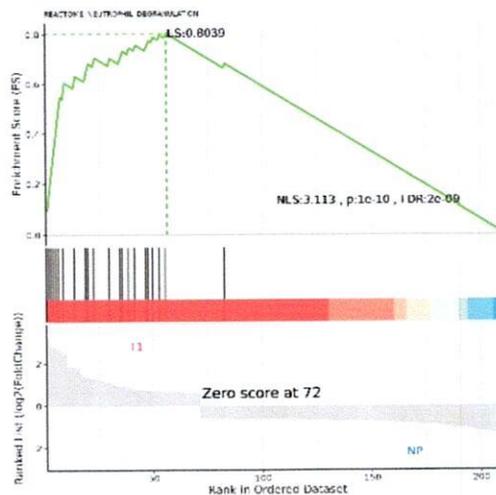
说明：横轴为比较组中的蛋白（蛋白集）根据其表达量变化值（Log2FC）由大到小的排序。GSEA 富集图自上而下分为三部分：①上部分显示的是当分析沿着蛋白集按排序计算时，ES（Enrichment Score）值在计算到每个蛋白位置时的展示(即分析过程中动态的 ES 值)，最高峰处的 ES 得分即为该通路的 ES 值。②中间部分俗称条形码图，用线条标记了该通路中涉及到的蛋白出现在蛋白集排序列表中的位置。红蓝相间的热图是表达丰度排列（红色越深的表示该位置的基因 log2FC 越大，蓝色越深表示 log2FC 越小）。③最下面部分为排序后比较组中蛋白 Log2FC 的排序，以灰色面积图展示。图的右上侧的注释为 KEGG 通路富集 pvalue 和 FDR 值。

输出文件：

- 1) 3-3-7 GSEA KEGG 分析

### 3.3.7.3 GSEA Reactome 富集分析

通过 GSEA 的方式将可定量蛋白进行 Reactome 通路富集分析。



### Groupvs 组 GSEA Reactome 富集分析

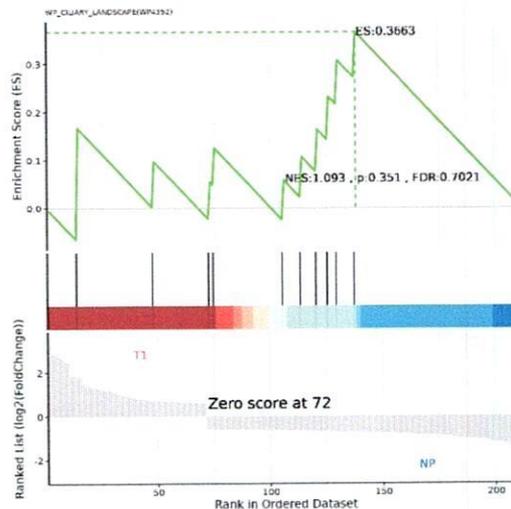
说明：横轴为比较组中的蛋白（蛋白集）根据其表达量变化值（Log2FC）由大到小的排序。GSEA 富集图自上而下分为三部分：①上部分显示的是当分析沿着蛋白集按排序计算时，ES（Enrichment Score）值在计算到每个蛋白位置时的展示(即分析过程中动态的 ES 值)，最高峰处的 ES 得分即为该通路的 ES 值。②中间部分俗称条形码图，用线条标记了该通路中涉及到的蛋白出现在蛋白集排序列表中的位置。红蓝相间的热图是表达丰度排列（红色越深的表示该位置的基因 log2FC 越大，蓝色越深表示 log2FC 越小）。③最下面部分为排序后比较组中蛋白 Log2FC 的排序，以灰色面积图展示。图的右上侧的注释为 Reactome 通路富集 pvalue 和 FDR 值。

输出文件：

#### 1) 3-3-7 GSEA Reactome 分析

### 3.3.7.4 GSEA WikiPathways 富集分析

通过 GSEA 的方式将可定量蛋白进行 WikiPathways 通路富集分析。



#### Groupvs 组 GSEA WikiPathways 富集分析

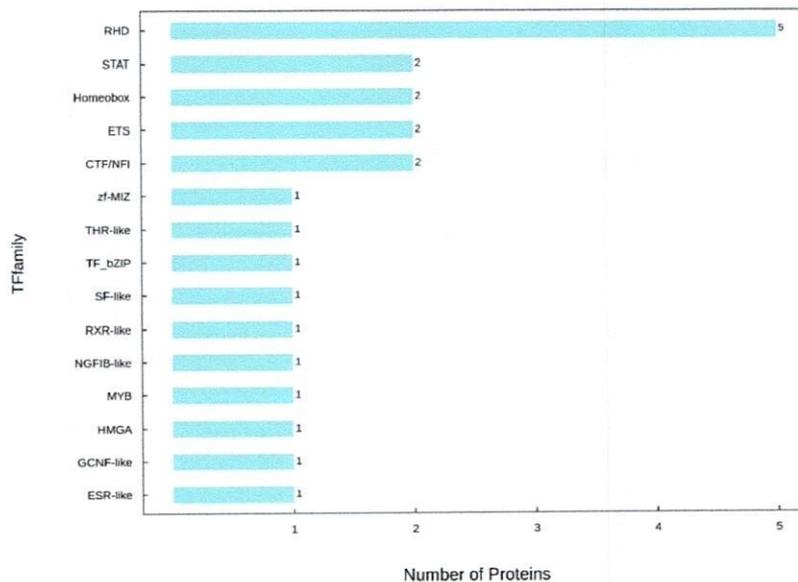
说明：横轴为比较组中的蛋白（蛋白集）根据其表达量变化值（Log2FC）由大到小的排序。GSEA 富集图自上而下分为三部分：①上部分显示的是当分析沿着蛋白集按排序计算时，ES（Enrichment Score）值在计算到每个蛋白位置时的展示(即分析过程中动态的 ES 值)，最高峰处的 ES 得分即为该通路的 ES 值。②中间部分俗称条形码图，用线条标记了该通路中涉及到的蛋白出现在蛋白集排序列表中的位置。红蓝相间的热图是表达丰度排列（红色越深的表示该位置的基因 log2FC 越大，蓝色越深表示 log2FC 越小）。③最下面部分为排序后比较组中蛋白 Log2FC 的排序，以灰色面积图展示。图的右上侧的注释为 WikiPathways 通路富集 pvalue 和 FDR 值。

输出文件：

#### 1) 3-3-7 GSEA WikiPathways 分析

### 3.3.8 转录因子分析

转录因子(Transcription Factor, TF)是指能够以序列特异性方式结合 DNA 并且调节转录的蛋白质, 由于转录因子有特殊的功能, 会对这类蛋白进行注释并进行深入分析。PlantTFDB (Plant Transcription Factor Database) 和 AnimalTFDB (Animal Transcription Factor Database) 数据库分别包含植物和动物的转录因子及转录因子家族信息, 可预测所关注的蛋白是否为转录因子, 以及所属的转录因子家族。



转录因子注释结果条形图

说明: 纵坐标代表转录因子家族, 横坐标代表注释到该转录因子家族的蛋白数目, 浅蓝色为注释到该转录因子家族的差异蛋白数量, 深蓝色为注释到该转录因子家族的鉴定所有蛋白数量。

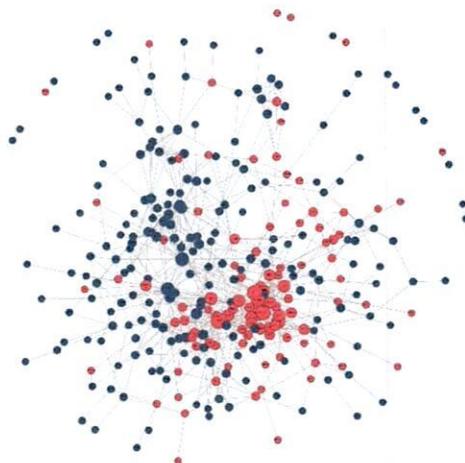
输出文件:

- 1) 3-3-8转录因子分析

### 3.3.9 蛋白互作网络分析

蛋白质发挥功能的重要方式之一就是与其他蛋白发生相互作用, 通过蛋白间介导的途径、或形成复合物进而发挥生物学调控作用。例如, 高度聚集的蛋白质可能具有相同或相似的功能; 连接度高的蛋白质可能是影响整个系统代谢或信号转导途径的关键点。因此研究蛋白-蛋白相互作用 (Protein Protein Interaction, PPI) 具有重要意义。此外, 将蛋白质相互作用网络分析和通路注释的结果相结合, 还可以获得更全面系统的分子层面的细胞活动模型, 便于分子机制的深入研究和挖掘。

本项目基于 STRING 数据库中的蛋白质相互作用关系，对 groupvs 比较组的差异表达蛋白质构建蛋白质互作网络图，如下图。



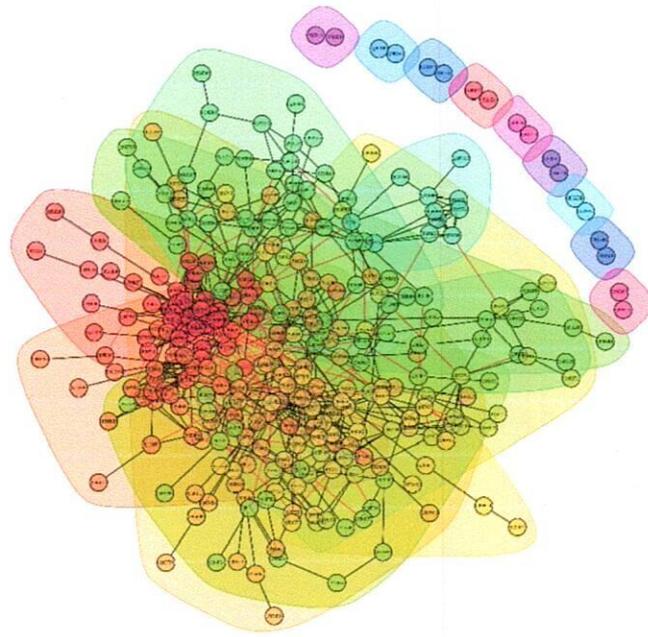
groupvs 组差异表达蛋白质相互作用网络

说明：图中圆圈结点表示差异表达蛋白质，线表示蛋白质与蛋白质之间的相互作用。其中圆圈颜色表示蛋白质表达差异（上调标注为红色、下调标注为蓝色、），圆圈大小表明该蛋白质连接度（即与某蛋白直接相互作用的蛋白质数目）。**通常来讲，连接度越大，该蛋白质发生变化时整个系统受到的扰动就越大，更可能是维持系统平衡和稳定的关键，为后续重点研究的候选蛋白质。**

输出文件：

#### 1) 3-3-9蛋白互作网络分析

在 PPI 互作网络中，高度聚集的蛋白质往往可能具有相同或相似的功能，并通过协同作用发挥生物学功能。因此，基于拓扑结构识别原理，将互作网络图中聚集程度高的蛋白划分为不同簇（Cluster）。具体划分簇展示如下图（每一类簇的展示图详见输出文件）。进一步，对每一类簇进行功能方向归类，功能归类表参见输出文件。



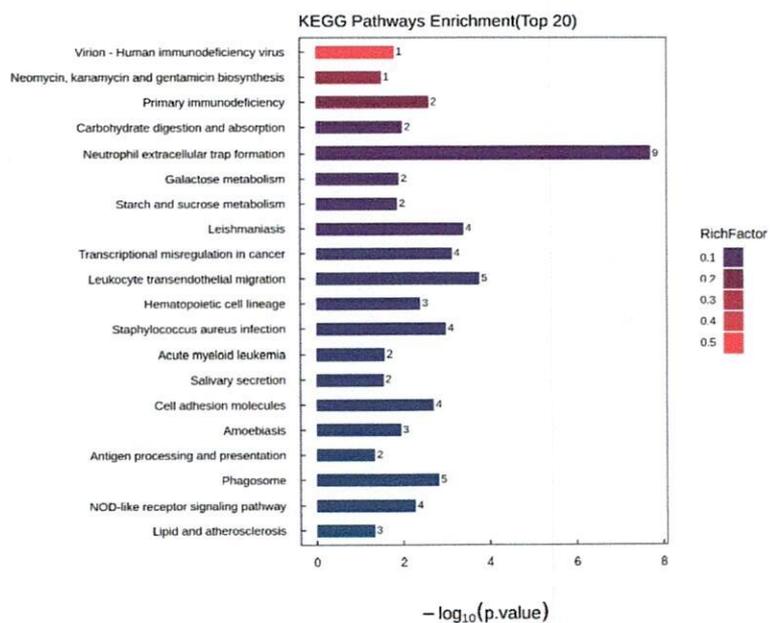
groupvs 组互作蛋白 module 分析图

说明：通常来讲，同一网络模块内蛋白往往具有相似的生物学功能，可选区感兴趣功能模块内的蛋白作为后续研究重点。

输出文件：

#### 1)3-3-9 蛋白互作网络分析

根据上图的高度聚集的蛋白质归类结果，从多到少分成多个 Cluster，选取前 5 个功能簇，每个 Cluster 中的相关高度聚集蛋白质进行 KEGG 富集分析，绘制 KEGG 富集条目柱形图和气泡图，文中仅展示一个 Cluster 的结果，其他结果见附件。进行 KEGG 通路注释统计，结果如下图所示。



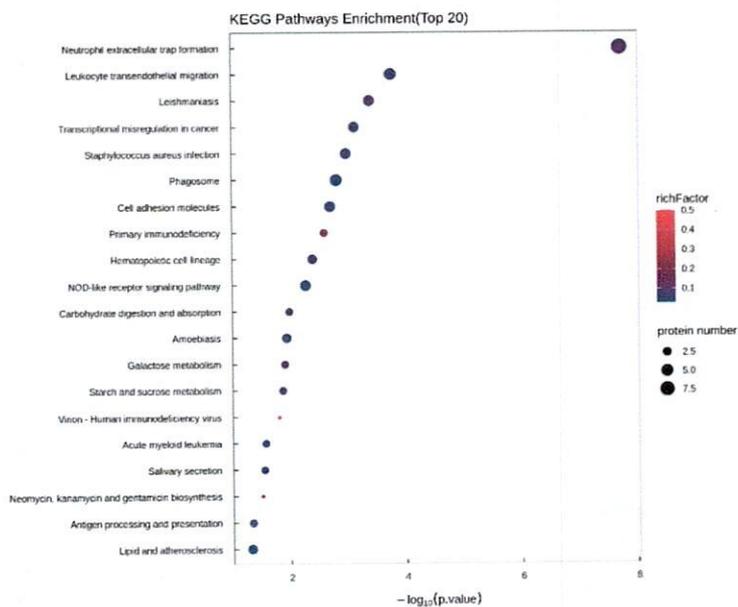
groups 组显著差异蛋白质的 KEGG 通路注释统计图 (Top20)

说明：图中纵坐标是含差异表达蛋白质参与的通路名称，横坐标表示富集显著性，即基于 Fisher 精确检验

(Fisher's Exact Test) 计算 P 值 (取 $-\log_{10}$ )，横坐标的值越大表示对应的通路下富集度的显著性水平越高。一般情况下，参与某一通路的差异表达蛋白质数目越多，说明该通路越重要，需要重点关注或者进行后续深入机制的探讨。

输出文件：

1) 3-3-9 蛋白互作网络分析



groups 组所有差异蛋白质的 KEGG 通路富集气泡图 (Top20)

输出文件：

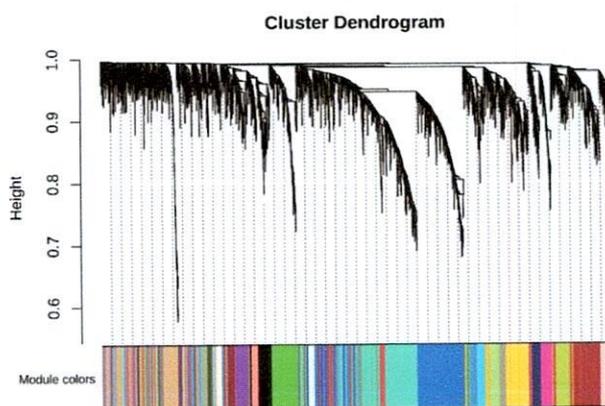
#### 1) 3-3-9 蛋白互作网络分析

### 3.3.10 表型与组学关联分析

WGCNA 分析 (weighted gene co-expression network analysis), 主要原理是通过加权共表达网络分析的方式分析多样本蛋白的表达模式, 鉴定出高度协同变化的蛋白模块 (module), 并根据蛋白模块的内连性和模块与特定性状或表型之间的关联, 筛选候选生物标记物或治疗靶点, 在疾病以及其他性状与蛋白关联分析等方面的研究中广泛应用。主要应用特色有两点: 一是将表型数据如性别、年龄或者其他生理特征与组学进行关联分析, 寻找表型与分子数据的关键的功能调控模块。二是通过基因间的表达相关性及权重挖掘关键功能分子, 识别已知基因的新功能。

#### 3.3.10.1 关键功能模块分析

首先, 基于最优软阈值计算蛋白间表达量相关系数, 构建蛋白层次聚类树划分共表达模块, 并对表达模式相近的模块进行合并, 获得最终划分的蛋白共表达模块:



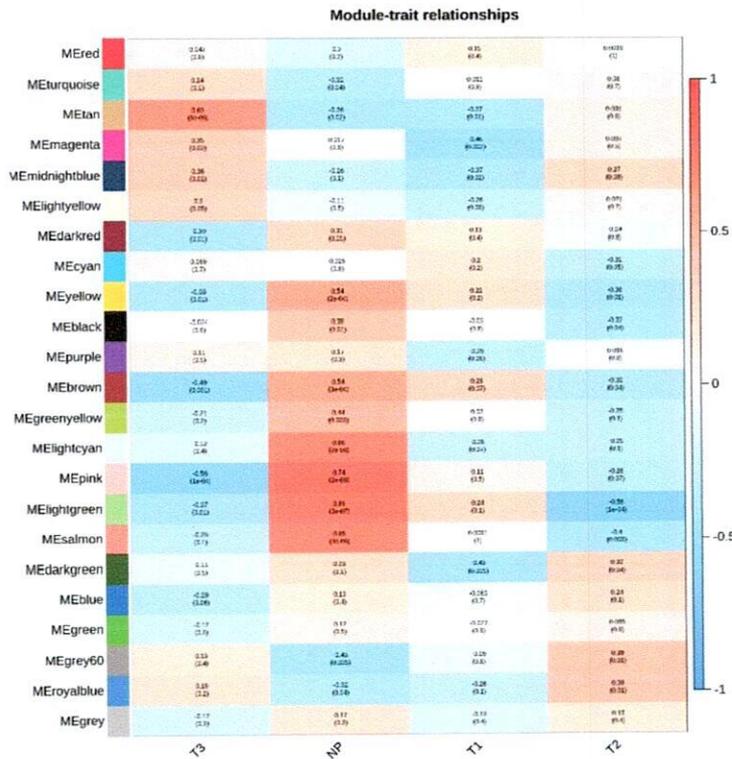
蛋白共表达模块划分图

说明: 横坐标展示的是不同的模块, 每个模块用不同的颜色表示, 灰色表示无法进行归类的模块; 每个树杈代表一个蛋白, 树杈的距离体现了蛋白的相似程度, 树杈越短相似性越高, 表达模式相近的蛋白聚集在一个分支里。

输出文件:

#### 1) 3-3-10 Module 构建

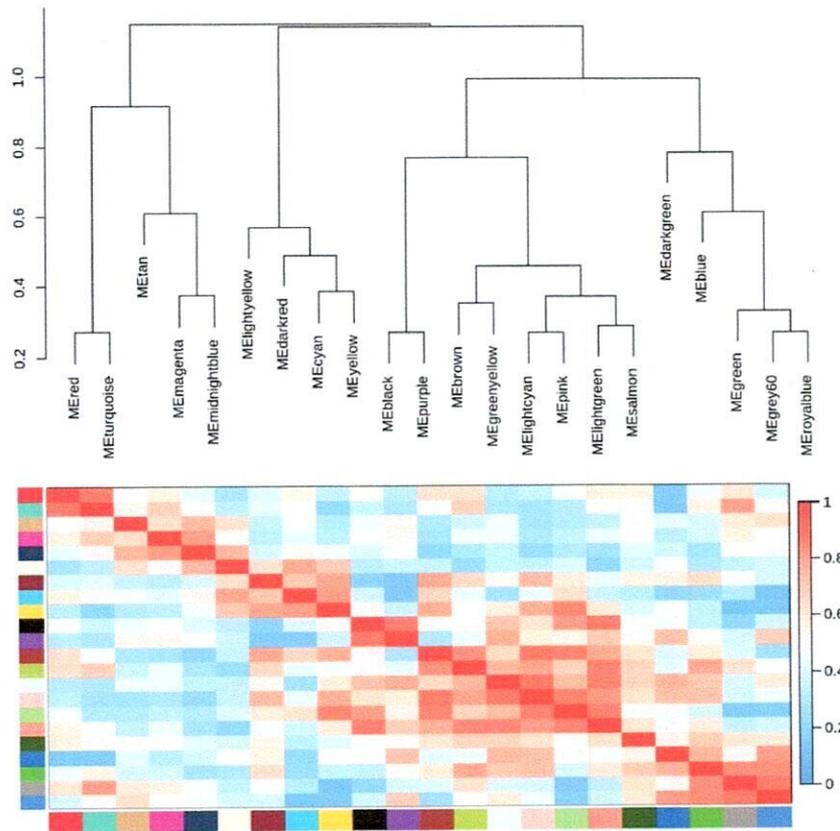
其次, 以分组或临床表型作为性状, 可以获得与该性状相关的共表达模块情况。进一步可从模块与表型性状相关性图中挑选感兴趣的性状所对应的模块, 一般筛选标准: 相关性数值越接近 $\pm 1$ , 相关性检验 P 值小于 0.05, 表明该模块是决定该性状的关键模块。



模块与表型性状相关性热图

说明：横坐标为性状或表型因素； 左边纵坐标的色块为不同的 module 类型（以不同颜色命名不同 module）； 中间的大色块代表各蛋白 module，同时在色块上标注了 p-value（括号中的值）及相关性系数值；右边 color bar 颜色代表相关系数大小，相关系数值介于-1 和+1 之间，红色代表正相关、蓝色代表负相关，相关性越高则颜色越深，相关性越低则颜色越浅。

再者，观察模块特征值的聚类树杈图，并通过聚类热图识别出表达模式更加相似的模块组。如下图所示：



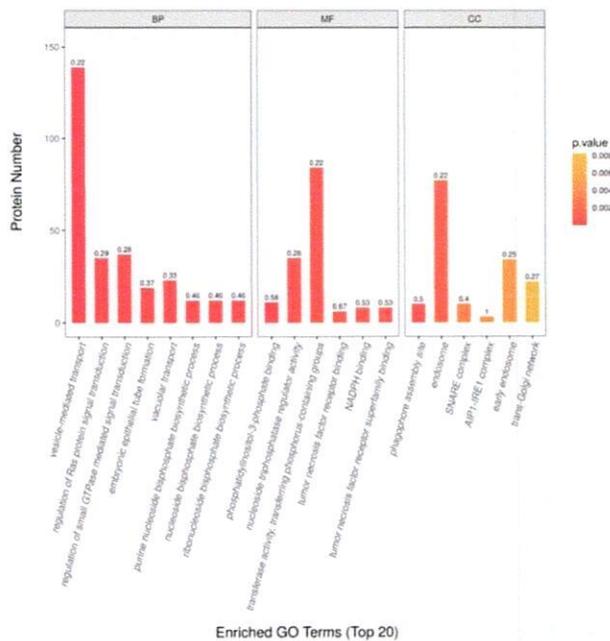
模块（性状）关联聚类图

说明：上半张聚类树杈图展示了不同模块（性状）之间的相异程度，距离越远则说明相异程度越高。下半张图展示了不同模块（性状）之间的相似程度，颜色越接近红色则说明相似程度越高

输出文件：

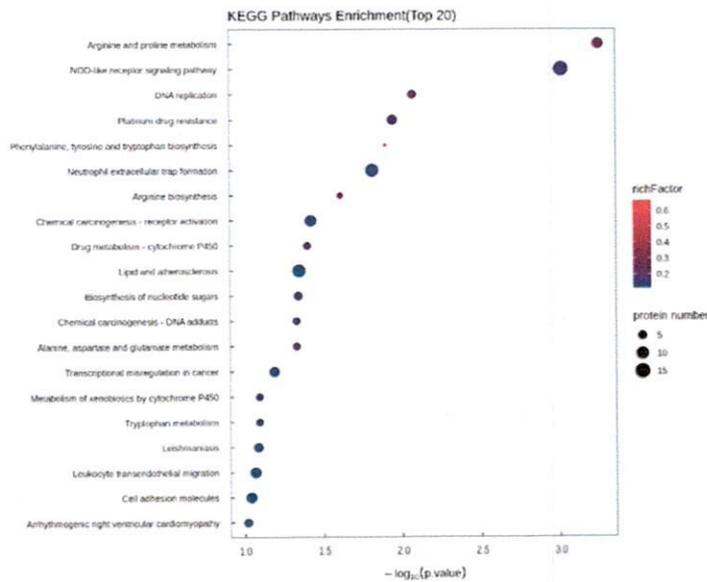
1) 3-3-10 Module 构建

将所挑选出的模块中的蛋白进行 GO/KEGG 富集分析，以查看该模块中共表达蛋白的功能/通路富集情况，反映该模块的蛋白功能/通路水平的特征，如下图所示：



某模块蛋白的 GO 富集分析图

说明：横坐标：GO 功能类别；纵坐标：与 GO 功能相关的蛋白质数量；条形图数字标签标识：富集因子 (Rich factor)；颜色深浅表示 p 值大小，即某个功能受到影响的显著程度。



某模块蛋白的 KEGG 富集分析图

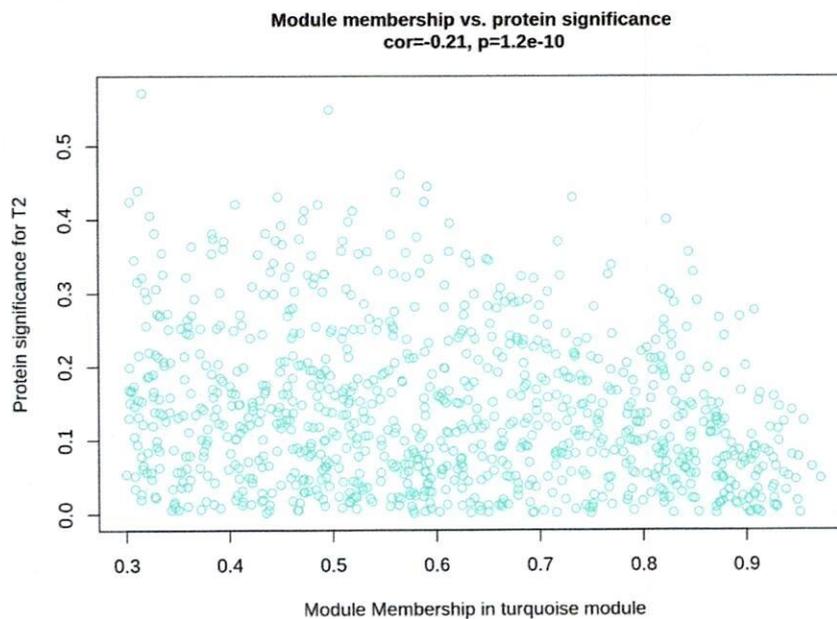
说明：图中横坐标为某 KEGG 通路的富集显著性，即基于 Fisher 精确检验 (Fisher's Exact Test) 计算 P 值 (取  $-\log_{10}$ )，横坐标的值越大表示对应代谢通路富集度的显著性水平越高，颜色梯度代表富集因子的大小 (Rich Factor $\leq 1$ )，富集因子表示注释到 KEGG 通路类别的显著差异表达蛋白质数目占注释到该类别的所有鉴定到的蛋白质数目的比例，颜色越接近红色代表 Rich Factor 值越大，气泡的大小表示每个 KEGG 通路下差异蛋白质数目。

输出文件：

### 1) 3-3-10 Module 构建

### 3.3.10.2 重要模块核心蛋白分析

在筛选出与性状（样本）高度相关的模块后，我们还可以观察 Gene Significance（蛋白与性状/样本的相关性）与 Module Membership（蛋白与模块的相关性）在每个模块中的散点分布情况，从而探究重点模块中蛋白与模块的相关性和基因与性状的相关性是否有较好的一致性，并且进一步从重点模块中筛选出相关性较高的 Hub protein（核心蛋白），如下图所示：



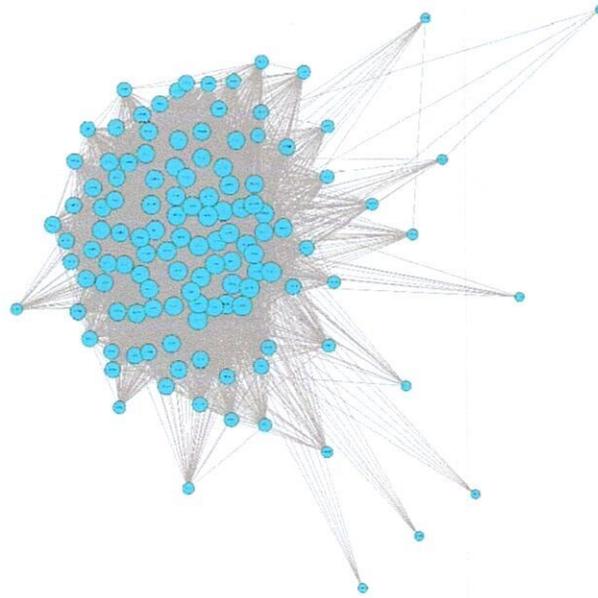
groupvs 散点分布图

说明：横坐标为每个蛋白的表达量与模块特征值的相关系数，纵坐标为每个蛋白的表达量与性状（样本）数据的相关性，cor 值为 GS 和 MM 值的相关系数，p 值则是对相关系数的假设检验值。

输出文件：

#### 1) 3-3-10 Module 构建

将所挑选出的模块中包含的所有蛋白构建蛋白共表达网络，进行可视化展示，反映模块中蛋白的相互关系，从另外一个角度筛选该模块共表达网络中的核心蛋白，进行后续深入研究。



### 模块内蛋白共表达网络分析

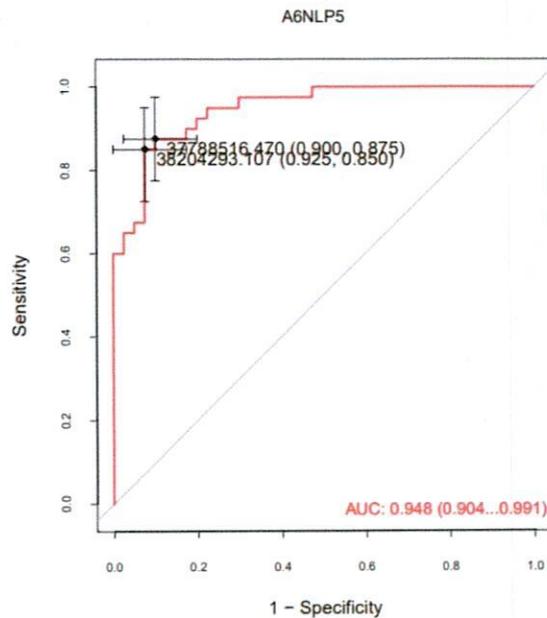
说明：圆圈颜色与模块颜色相同，如该图圆圈为蓝色，表示该互动网络为蓝色模块的蛋白构建的互动网络，节点大小与其连接度（Degree）正相关，即与该蛋白相关的蛋白越多，则其连接度越大，节点尺寸就越大，其在网络中的地位越关键。线条表示蛋白之间的互动关系，线条的粗细与相关系数的绝对值成正比，即线条越粗，相关程度越高。

输出文件：

1) 3-3-10 Module 构建

### 3.3.11 ROC 分析

ROC 分析（receiver operating characteristic curve，受试者工作特征曲线）是把灵敏度和特异度结合起来综合评价诊断准确度或判别效果的一种方法，在医学领域中广泛用于临床诊断、人群筛选等研究。在蛋白质组学中，对比较组间的差异蛋白进行 ROC 分析，可以用来指示 biomarker 区分两组间（实验组和对照组）的能力，展示 AUC 值 TOP25 的蛋白。



### ROC 分析

说明：横坐标为 1-特异性，即假阳性率——阴性群体中，检测为阳性的概率，希望该值越低越好；纵坐标为敏感度，即真阳性率——阳性群体中，检测出阳性的概率，希望该值越高越好；曲线越往左上角说明预测准确率越高，曲线下面积越大，即 AUC 值越大说明预测准确率越高。图形中红色标识的文字为该曲线对应的 AUC 值和 95%的置信区间；黑色标识的文字为原始的强度表最佳临界值，括号中位特异度和灵敏度。

输出文件：

1) 3-3-11 ROC

## 附件三 售后服务

### 1 服务内容

- 1.1 提供专业的学术顾问咨询给出文章投稿建议，并提供文章杂志要求的分辨率和大小合格的图表；
- 1.2 对审稿人的修回意见提供专业答复供参考；
- 1.3 定制化信息分析：完成项目方案内的定制化信息分析内容
- 1.4 定制化绘图：在项目合作期内，为合作伙伴发表项目产出论文提供方案内的高端定制绘图服务

1.5 应急解决方案：公司设立技术支持领导小组保证突发事件发生时,能够迅速召集技术人员,立即制定应急技术方案对一般性技术故障,如果检测意外失败、实验意外失败（未达到实验结果的质控标准），可以评估后安排立即复测。不另计费用。

1.6 其他售后服务：思路拓展、投稿建议、生信云平台等。

## **2 培训与维护**

2.1 根据需求提供技术培训及咨询服务；对 DIA 蛋白组学进行培训。

2.2 技术咨询与指导：为招标方项目参与人员提供技术咨询与指导，包括质谱原理，实验操作，数据处理以及其他与本项目相关的技术指导

2.3 提供至少 1 次组学分析培训交流

## **3 售后服务方式**

3.1 电话服务，能 7\*24 小时处理客户问题，全国范围内有相应的销售与技术，免费进行疑难问题解答。

3.2 远程连接技术支持人员，通过腾讯会议等在线会议方式，对数据分析结果进行指导。

3.3 往来信函、传真、电子邮件，解答用户在使用中碰到的各种技术问题。

3.4 现场服务：在客户授权的情况下，针对客户已有数据结果进行分析，提供解决方案。

3.5 定期汇报：项目进度以及数据分析处理进程。

## **4.服务响应时间**

4.1 我们将对用户 provide 全方位的售后服务,并提供最佳的服务响应时间。

4.2 电话服务技术支持与服务时间为 8:30-17:00,周一至周五(国家法定的休息日和节假

4.3 数据验收合格后，至少提供 2 年的售后服务；

4.4 7\*24 小时处理客户问题

## **5.保密要求**

未经甲方同意，乙方不得以任何形式向甲方指定的负责人以外的单位或者个人披露项目内任何数据情况，并且有义务协助保存相关数据。

---

如果一方因另一方违反其在本协议下的义务而使披露方遭受直接损失，包括但不限于财产毁损、付费和支出以及诉讼费用（包括但不限于律师费用）时，违约方应对守约方予以补偿，以使守约方免受其违反保密承诺行为的损害。

附件四 中标通知书



北京科技园拍卖招标有限公司  
Beijing Science Park Auction & Tender Co., Ltd

中标通知书

上海中科新生命生物科技有限公司

根据【项目编号：KJY20250666】基于人工智能和健康大数据的急性心血管事件预测和预警研究项目（第1包）的招标文件和你单位于2025年5月22日提交的投标文件，经评标委员会综合评审，最终确定你单位为中标单位，中标金额人民币1100000元，大写：壹佰壹拾万元整。

请你单位在本中标通知书发出之日起30日内，按照采购文件确定的事项与采购人签订采购合同。

在你单位与采购人签订采购合同后5个工作日内，携带本通知书和贵公司开出的退投标保证金收据与我公司联系办理退还投标保证金事宜。

代理机构：北京科技园拍卖招标有限公司（盖章）



日期：2025年5月28日

地址：北京市海淀区万泉庄万柳光大西园6#楼0188邮编100089  
Tel: 010-82575731 Fax: 010-82575350/5840 Http://www.bkpmzb.com