

# 2021 年腾讯微信犀牛鸟专项研究计划课题列表

## 目录

1.	社交网络与图计算 .....	2
1.1.	个人知识图谱的表示与学习 .....	2
1.2.	基于图神经网络的推荐系统算法研究与应用 .....	3
1.3.	稀疏空间数据挖掘 .....	3
1.4.	欺诈团伙挖掘 .....	4
2.	短视频及关联技术 .....	4
2.1.	视频号智能红点的强化学习混排 .....	4
2.2.	视频号推荐算法优化 .....	5
2.3.	短视频推荐的多目标学习研究 .....	5
2.4.	基于因果推理的推荐方法探究 .....	6
2.5.	推荐场景中结合视频多模态&用户行为的主题模型 .....	6
2.6.	消费级视频元学习对抗鲁棒性研究 .....	7
2.7.	多模态视频检索任务 .....	7
2.8.	短视频危险动作识别 .....	8
3.	支付安全与风险防空 .....	8
3.1.	联邦学习鲁棒性与多方计算平台模型研发 .....	8
3.2.	金融反欺诈领域的大规模图计算研究 .....	9
3.3.	信贷风控算法研究 .....	10
4.	搜索与推荐 .....	10
4.1.	基于自监督学习的预训练搜索排序模型 .....	10
4.2.	基于用户行为的搜索个性化排序建模技术研究 .....	11
4.3.	结合领域知识的机器学习算法研究 .....	12
5.	自然语言处理 .....	12
5.1.	自然语言处理和对话系统前沿技术研发与应用 .....	12
5.2.	从文本到 SQL (Text2SQL) 的自动语义解析 (Semantic Parsing) 算法研究 .....	13
6.	实验平台与软件测试 .....	13
6.1.	实验平台算法模型研发 .....	13
6.2.	基于微服务依赖图的故障根因定位 .....	14

# 1. 社交网络与图计算

## 1.1. 个人知识图谱的表示与学习

近年来，随着深度学习在 NLP 领域的突破，知识库构建、实体链接、智能问答等技术得到了进一步提高。然而，当前关于知识图谱的研究往往集中于全局的实体和关系，较少建模用户的个人兴趣及知识。谷歌于 2019 年提出了“个人知识图谱 Personal Knowledge Graphs”的概念，对用户的个性化知识表示、存储、建模提出了挑战。我们提出如下科研问题：

- ◇ 如何利用全局知识图谱及用户行为，有效构建用户的个人知识图谱？
- ◇ 如何对个人知识图谱进行有效的表示、学习、应用？
- ◇ 如何建模用户的兴趣，包括判别长期与短期兴趣、已知与未知知识等？

具体研究课题包括但不限于以下方面：

- ◇ 个人知识图谱的表示与学习；
- ◇ 个人知识图谱中的用户兴趣建模、长短期兴趣判别、已知与未知知识的表示；
- ◇ 垂直领域，私域实体和关系的少样本、零样本抽取；
- ◇ 基于个人知识图谱和行为序列，通过注意力机制、预训练等方法，进行用户行为预测；
- ◇ 基于个人知识图谱改进搜索和推荐系统、打破“信息茧房”，提升信息多样性等。

### 科研目标：

基于学术界和业界的最新研究，探索个人知识图谱表示与学习的前沿技术，并在我们提供的匿名化数据集或学界的开放数据集上进行实验。

产出 NLP 领域顶级会议论文 1 篇，技术落地或部分落地到生产环境，并支持成果发表专利。

### 可提供资源：

匿名化训练及测试数据集，已有的全局知识图谱和计算资源，技术辅导。

**关键词：**机器学习，自然语言处理，知识图谱。

## 1.2. 基于图神经网络的推荐系统算法研究与应用

推荐系统架构越来越复杂，输入数据结构也越来越多样，其中包括了多源异构兴趣网络 and 用户社交网络等图结构数据，之前的方式都是通过将图数据转化为欧氏空间结构数据，并用深度学习方法处理。图神经网络（graph neural network, GNN）具有直接从图的领域对数据进行特征提取和表示的优势，并在一些领域证明了其可行性。所以，将图神经网络和深度学习结合产出一个端到端的推荐系统，是一项有难度并且有意义的工作，也是推荐系统新的研究热点。

### 科研目标：

提出新型 GNN 和与推荐系统相结合的方案，在推荐场景下融入用户的社交网络和多源异构网络信息，实现大规模并行计算算法，并且在短视频推荐业务上落地，相对提升线上效果不少于 5%，产出顶会论文。

### 可提供资源：

- 1、脱敏的短视频业务用户历史行为数据集；
- 2、脱敏的多源异构网络和社交网络的数据集；
- 3、计算资源和图计算平台。

**关键词：**推荐系统，图算法，GNN。

## 1.3. 稀疏空间数据挖掘

微信支付在广告、营销、支付分等应用中，对商户画像和城市地块特征挖掘有强烈需求。学界中，通过签到数据和GIS数据集挖掘财富指标和地块发展指标的方案已被众多研究者证实是可行的，并发表于Science、Nature Communications等顶级刊物上。然而，由于轨迹数据稀疏性、GIS数据集不完备等限制因素（如：缺少遥感影像、城市地价信息不全），无法直接应用现有的方案，还需要投入大量时间和人力进行数据收集和算法改进。本项目旨在借助高校的现有研究成果，实现如下快速突破：

- 1、对空间数据融合和预测的方案，并用于生产地块尺度的商业活力指标、房价和欺诈热点等；
- 2、基于稀疏轨迹embedding和时空上下文感知，用于评估资产和风险等特征。

### 科研目标：

提出一种地块信息融合和预测机器学习模型，并生产全国的地块特征；提出一种稀疏轨迹embedding和时空感知方案，用于挖掘资产和风险特征。

产出 2 篇一区期刊论文（JCR分区）和 2 个专利。

**可提供资源：**

可提供部分数据、服务器资源和驻场实习生职位。

**关键词：**LBS，时空数据挖掘。

## 1.4. 欺诈团伙挖掘

微信平台中存在欺诈分子，通过各类手段对普通用户进行诈骗，这会影晌正常用户的体验与平台健康。欺诈分子在作案时，会在网络上留下蛛丝马迹，在某些途径上形成关联关系。本课题将开展相关的前沿研究和工程创新。

**科研目标：**

本课题主要目标为技术落地，包括：

- 1、异构数据源的社团划分，同时结合考虑好友、群、IP等数据源对欺诈分子进行团伙聚合；
- 2、可解释的欺诈用户识别或欺诈交易识别，综合考虑各种图数据源的信息，提供高准确率和可解释的模型识别方案；
- 3、欺诈关键节点、关键路径和重要子图识别，分析欺诈节点的拓扑特征；
- 4、可视化，微信存在大量关联路径，二度扩散后常会超过十万级别的节点，如何快速定位可疑节点并提供知识发现的能力；
- 5、鼓励在顶级学术会议发表论文，以及专利产出。

**可提供资源：**

- 1、实时图数据库集群资源；
- 2、离线计算集群资源以及分布式图计算基础框架。

**关键词：**图计算，安全风控。

## 2. 短视频及关联技术

### 2.1. 视频号智能红点的强化学习混排

视频号智能红点聚焦于如何在减少用户打扰的前提下，提升对视频号活跃和消费的转化。智能红点涉及到好友点赞、关注、消息、关注直播、好友观看直播等多路红点的混合排序，产生方式、触发时机和业务目标都各有不同，候选集中的分布也有较大差异，很难通过传统机器学习的排序思路来解决。强化学习可以补全传统机器学习所不具备的

探索能力和最大化长期收益的能力，预期可以带来红点整体目标的进一步提升。

**科研目标：**

落地视频号智能红点的混排算法，提升视频号活跃占比 1%以上。

**关键词：**机器学习，深度学习，强化学习。

## 2.2. 视频号推荐算法优化

微信视频号推荐目前仍有较大的优化空间。例如，可以利用在线用户反馈来快速迭代模型，利用微信关系链数据来进行人群协同优化效果，利用微信朋友圈等数据进行画像挖掘等，我们希望在微信视频号推荐这个特定情境下做出创新。

**科研目标：**

在微信视频号推荐系统中落地应用，预期提高用户留存、时长等产品指标，并产出顶会论文两篇。

**可提供资源：**

脱敏的视频号用户数据，基线系统与服务器资源支持。

**关键词：**推荐系统。

## 2.3. 短视频推荐的多目标学习研究

在短视频推荐业务中，为了更好的兼顾用户体验与内容生产者的收益，往往需要优化多个不同方向的业务指标。我们认为多目标学习方法在实际应用中仍存在以下难题：

1、网络结构的设计问题：需要平衡各目标之间的学习稳定性、子模型表达能力；

2、样本中存在的标签不均衡及偏差 (bias) 问题：实际业务中不同目标的样本标签不均衡，且受到不同 bias 的影响，在多目标学习场景下需要有效解耦；

3、考虑任务难易程度的自监督学习：不同目标之间天然存在学习难易程度的差异，需要从自监督学习的角度，在不借助额外标签数据的条件下设计更加合理的学习范式。

**科研目标：**

提出新型的短视频推荐多目标学习方案，针对多目标学习网络结构的设计问题、样本中存在的标签不均衡及偏差 (bias) 问题、考虑任务难易程度的自监督学习问题等展开研究。该方案能够在短视频推荐业务中取得STOA的性能，对比现有业界解决方案（如MMoE）在用户互动率、停留时长、留存指标上有提升。

**可提供资源：**

脱敏的短视频业务用户历史数据，模型离线训练服务器资源，线上实验平台环境。

## 2.4. 基于因果推理的推荐方法探究

目前，基于机器学习、深度学习的推荐模型已经成为了推荐系统的主流方法，这类方法旨在更好地拟合用户行为数据，从而挖掘用户可能感兴趣的物品。然而，推荐系统收集到的用户行为数据往往是有偏差的，会受到用户的选择、视频的时长、商品的流行度等影响，因此用户的行为数据往往不能够真实反应出用户的偏好，这类基于数据驱动的方法难以取得预期效果。

因果推理旨在挖掘实体之间的因果关系、预测行为的因果效应并对观测的现象给与解释，是近年来的研究热点。将因果技术应用于推荐系统中，可以有效地挖掘并减缓推荐系统中的数据偏差问题。

### 科研目标：

本项目成果主要应用在视频号业务中，解决推荐场景中的bias问题，提升推荐点击率和转换率。

### 可提供资源：

经过匿名化处理的视频号用户行为数据、机器学习平台。

**关键词：**推荐系统。

## 2.5. 推荐场景中结合视频多模态&用户行为的主题模型

在推荐领域中，对视频向量一般使用ID embedding来表示，目前主流对embedding的表征分为两种，一种是通过用户行为（如曝光点击等行为）的 supervised learning，通过行为预训练用户的表示，我们称为动态表征。一种是通过视频进行标注&分类，通过分类模型得到视频的embedding表示，我们称为静态表征或者side information。是否能扩充一种自监督或无监督的方案，通过结合视频的多模态信息和用户行为，对视频进行聚类 and 表征，对视频生成隐式主题表征，该表征具有一定的解释性和受众匹配的效果，不同的主题能对应到不同的受众，是本课题的研究重点。

### 科研目标：

研究产出推荐场景中结合视频多模态与用户行为的主题模型，减轻视频冷启动问题，并提高推荐的准确率；同时由于主题具有一定的解释性，也可以作为召回逻辑应用到相关视频等业务场景中。

研究成果期望具有学术突破和行业创新性，发表顶会论文至少一篇。

### 可提供资源：

微信公众平台脱敏业务数据，Spark & GPU计算平台。

**关键词：**推荐系统。

## 2.6. 消费级视频元学习对抗鲁棒性研究

随着视频推荐业务的高速发展，视频内容的结构化和元学习已经被广泛应用到视频制作、审核、分发、消费等各个环节。然而在内容推荐任务当中，消费级视频的元学习面临着鲁棒性和可信性不足的问题。由于推荐场景下视频具有数据量大、场景复杂、类别多样的特点，原有常规的深度学习模型会呈现两个短板：模型鲁棒性和可迁移性不足，当推荐场景/内容/目标变化时，元学习表示的效果会大幅下降；模型可解释性不足，整体被当做黑箱使用。

因此，在实际应用场景中，在元学习中引入噪声样本，基于强多样性对抗攻击模型，鼓励网络决策边界的差异化，可以显著提高模型的可迁移能力，减少高层数据噪声，解决看一看众多推荐场景下的视频元学习鲁棒性问题。

**科研目标：**

本课题拟从应用与数据的角度出发，针对当前视频元学习面对多应用场景下复杂数据鲁棒性不足的核心难点，提出引入强多样性的对抗攻击模型，建立元学习鲁棒性和可解释性评测的指标体系，为面向信息推荐的视频元学习高效优化和可迁移能力提供理论和技术保障。主要指标包括：

1、在国际顶级会议或期刊上发表论文 1-2 篇；

2、实现基于对抗攻击的视频元学习算法，解决推荐系统中动态多样场景下模型鲁棒性的问题，并争取在线上产品获得应用。

**关键词：**推荐系统。

## 2.7. 多模态视频检索任务

随着互联网和移动设备的发展，网络上视频内容越来越多，用户浏览视频内容的时间越来越长。面对海量视频内容，文本-视频搜索（text-video retrieval）是用户访问视频的主要方式之一。本课题主要研究如何合理利用文本、视频以及点击数据，从多个维度融合，提升视频搜索的体验。

**科研目标：**

1、文本-视频检索任务中，query是文本，检索目标是视频，研究建立跨模态数据的相关性的相关方法；

2、已有方法依赖于带有caption/title的视频数据集作为监督语料来建立模态间的

联系，但在很多实际业务场景中（视频号等UGC短视频平台），视频没有语义清晰的caption或者title。研究利用搜索点击行为的弱监督信息来进行学习的方法；

3、在多模态检索场景下，研究如何使用模态间通用的知识图谱/实体/标签等外部信息来辅助搜索排序。

**可提供资源：**

可提供脱敏的数据资源：

1、视频数据总量：3亿；

2、视频静态属性：粉丝、脱敏作者信息等；

3、行为数据：用户与账号的关注关系、用户播放记录、点赞记录、收藏记录、打赏记录、评论记录等。

**关键词：**视频搜索。

## 2.8. 短视频危险动作识别

短视频中经常出现一些跑酷、赛车或者极限运动等，在没有安全措施的情况下，观众容易盲目模仿而受伤。因此，需要针对此类危险视频特征，给出高效的自动审核识别方法，从而指导微信内容审核团队给出“请勿模仿”的官方警告。自动发现危险动作视频，可以大大节省审核人力，提升用户体验。

**科研目标：**

本课题的目标偏落地实施，期待应用于视频号的自动审核机制中，对于危险动作视频，能够协助微信官方自动给出告警，对危险动作视频的识别准确率90%以上、召回率90%以上。

**关键词：**计算机视觉。

## 3. 支付安全与风险防空

### 3.1. 联邦学习鲁棒性与多方计算平台模型研发

由于联邦学习具有分布式的特性，模型会在多个私有计算环境上进行训练，并且这些环境的数据私有，不可被检查。目前联邦学习在鲁棒性问题上还存在诸多缺陷，实际应用还会面临更多复杂且隐蔽的攻击问题，比如定向攻击和非定向攻击的攻防方法，提高联邦学习的鲁棒性至关重要。因此，本课题的第一个研究内容为联邦学习中存在的攻击以及相应的防御方法。

另一方面，多方计算平台模型训练优化，隐私和效率通常是首要考虑因素，研究加



速模型训练的方法，在有限计算资源下如何增强特定模型的训练稳定性和训练性能，提升联邦学习效率，比如分布式模型优化、多方训练协议优化等，是本课题的另一个研究内容。

**科研目标：**

综合运用分布式、数据安全等方法，提升联邦学习多方计算平台的稳定性和鲁棒性，产出至少一篇学术论文、并在生产环境中应用。

**可提供资源：**

测试数据集和计算资源。

**关键词：**统计，计算机，安全。

### 3.2. 金融反欺诈领域的大规模图计算研究

微信是一个日活跃用户超过 10 亿的社交平台，用户每天都会产生很多行为数据，包括文本、图片、视频、表情等丰富的多模态数据；同时，微信用户覆盖人群广泛，无论年龄或性别，基本涵盖各行各业人群。由于每个人的知识背景和职业属性不同，不同的人对搜索认知和敏感度也不一样，在搜索场景下很容易划分为不同的圈层人群。微信搜一搜业务产品主要解决在微信生态场景下用户查找信息效率问题。在传统搜索引擎里面，主要是对用户查询和文档之间的相关度匹配程度进行建模。为了充分发挥微信数据特色和生态价值，提升不同圈层用户搜索体验，如何将不同圈层用户信息融入排序算法中进行建模是一个值得研究的课题。

在微信支付的业务场景，反欺诈是非常重要的主题，无论是针对个人的风险识别，还是针对小微商户的风险评估，图计算和关系网络异常识别都扮演着重要角色。个人/小微商户数据异构，维度多而稀疏，且难以被获取，而目前的风控方法都基于个人/小微商户自身的画像特征。通过结合微信用户之间的关系可以作为用户信息的补充，建立高效的图神经网络算法去学习关系数据表达，作为规则策略和有监督学习的互补方案，补足这些方法的短板，是本科研项目所解决的核心问题。

**科研目标：**

- 1、构建微信支付的人-人关联、人-物关联，基于时空的异构关系网络；
- 2、研究大规模图计算算法，识别关系网络中的异常风险；
- 3、研究图计算在反欺诈场景的应用，如洗钱，套现，诈骗等。

**可提供资源：**

可提供部分脱敏数据、服务器资源和驻场实习生职位，数据与资源要求以实习生的形式在公司内部使用。

**关键词：**图计算，复杂网络，反欺诈，异常识别，知识图谱。

### 3.3. 信贷风控算法研究

随着微信支付金融业务布局的逐步完善，通过算法及模型降低信用及欺诈风险，在提升用户体验和产品收益中扮演重要角色。

近年来，信贷风控的研究逐渐出现在顶级学术会议的视野中，然而由于风控算法中使用的假设通常与业务数据特点和风险行为具有较强耦合性，导致许多前沿研究在实际应用中效果不尽如人意。例如金融信贷风险具有隐蔽性（如套现、养号、刷单等场景，无法获取准确无偏的样本标签）、对抗性（风险用户会改变行为试探攻击模型，导致模型失效）、延迟性（风险表现存在较长观察周期）等问题，这些问题限制了算法的实际表现。

本课题旨在结合微信支付自身金融业务特征和风控痛点特征，通过高校合作设计可落地算法，提升以下信贷领域的风控算法能力：

1、信用风险预测：优化信用模型算法框架，解决样本迁移、延迟性问题，提升评分卡、额度、息费、响应、还款能力、还款意愿模型的应用效果；

2、信贷欺诈识别：针对套现、养号、刷单、恶意透支等信贷风控领域的欺诈问题，设计高对抗性算法，识别欺诈行为并实时防控。

#### **科研目标：**

结合微信支付具体金融风控问题，包括但不限于：信用推断、套现识别、刷单养号识别、信贷欺诈识别等；降低整体信贷违约风险 20%，降低欺诈发生率 20%。

知识产权方面，计划共同发表 2 篇一区论文，产出 2 个专利。

#### **可提供资源：**

可提供部分脱敏数据、服务器资源和驻场实习生职位。

**关键词：**信贷欺诈识别、信用风险预测。

## 4. 搜索与推荐

### 4.1. 基于自监督学习的预训练搜索排序模型

在搜索排序模型建模过程中，需要对排序目标进行标注，但由于标注成本高，业界一般采用的方法是利用用户点击行为日志来快速、自动地构建标签数据。然而用户点击行为是有偏差的，并且是稀疏长尾的，直接利用可能造成负面影响。得益于自监督学习任务的发展，图像和自然语言处理领域都取得了具有突破性的进展。自监督学习逐渐成

为备受关注的应对标注数据缺乏的热门解决方案。如何从大规模的、带噪音的用户行为日志中构建与排序学习相关的自监督任务，学习好query、doc的表征，最终提高排序性能是一个值得研究的科研课题。

**科研目标：**

探索自监督学习在LTR领域的创新应用及相关的技术储备，具体产出包括：

- 1、业务效果的显著提升；
- 2、合著学术论文 1~2 篇，并投稿至相关的顶级学术会议（CCF A或B类会议，如SIGIR、WWW、CIKM、WSDM等）、1-2 篇专利；
- 3、形成模型及算法工具包一份。

**可提供资源：**

- 1、脱敏用户搜索日志；
- 2、Case debug工具；
- 3、训练模型所需的软硬件资源，如CPU集群、GPU集群。

**关键词：**搜索排序。

## 4.2. 基于用户行为的搜索个性化排序建模技术研究

微信是一个日活跃用户超过 10 亿的社交平台，用户每天都会产生很多行为数据，包括文本、图片、视频、表情等丰富的多模态数据；同时，微信用户覆盖人群广泛，无论年龄或性别，基本涵盖各行各业人群。由于每个人的知识背景和职业属性不同，不同的人对搜索认知和敏感度也不一样，在搜索场景下很容易划分为不同的圈层人群。

微信搜一搜业务产品主要解决在微信生态场景下用户查找信息效率问题。在传统搜索引擎里面，主要是对用户查询和文档之间的相关度匹配程度进行建模。为了充分发挥微信数据特色和生态价值，提升不同圈层用户搜索体验，如何将不同圈层用户信息融入排序算法中进行建模是一个值得研究的课题。

**科研目标：**

- 1、围绕基于用户行为的个性化搜索技术的探索，能输出论文至少 2 篇；
- 2、该项目研究成果主要应用并落地于微信搜一搜产品中。

**可提供资源：**

数据侧：搜一搜点击日志数据、微信圈层用户人群画像数据，因为涉及到用户数据隐私，所有数据均脱敏（建议研究合作者到公司内部集群开发环境进行模型训练和效果评测）；

工具及机器资源：大规模分布式（深度学习）训练平台、大数据处理分析平台等。

**关键词：**用户行为分析与建模，个性化排序技术。

### 4.3. 结合领域知识的机器学习算法研究

在微信支付的推荐、营销、权益、积分商城等应用中，机器学习扮演着重要角色。传统的机器学习算法往往需要比较丰富的样本和特征数据才能有较好的效果，而实际应用中这个条件并不总能满足。比如在行业推荐中，新商户接入时提供的数据有限；新产品拉新时，可用于识别高潜用户的样本也有限。另一方面，产品、运营和研发人员有一定的人工经验并可通过其他渠道学习获得一定的知识。如何把这些经验和知识更好的引入到机器学习建模的过程中，达到举一反三的效果是本课题所要解决的核心问题。

**本项目研究课题：**

- 1、如何把人工经验和分析结论组织成领域知识；
- 2、如何把领域知识引入到机器学习并提升效率和效果；
- 3、如何通过领域知识更好的解释和运用机器学习模型。

**科研目标：**

本项目希望通过科研合作更好的解决微信支付实际业务中机器学习应用效果、效率和可解释性的问题，提升微信支付推荐、营销、权益和积分商城中的推荐营销效果。

知识产权方面，期待产出 2 篇顶会论文和 2 个专利。

**可提供资源：**

可提供脱敏数据，以及服务器资源等。

**关键词：**结合领域知识的机器学习算法，机器学习，领域知识，推荐算法。

## 5. 自然语言处理

### 5.1. 自然语言处理和对话系统前沿技术研发与应用

本课题致力于研发自然语言处理和对话方向的前沿技术，用于实际业务场景，同时探索从感知到认知的新一代 AI 形态，占得未来技术先机，具体涵盖如下主题：

- ◇ 自然语言处理中基础技术前沿算法研发，包括但不限于：通用预训练技术，序列标注模型、文本分类、文本匹配、文本摘要、文本生成等先进模型研发；
- ◇ 对话系统中的前沿算法技术研发，包括但不限于：语义理解、状态跟踪、回复生成等前沿技术研发；
- ◇ 探索 AI 认知新方法和语料资源构建，包括但不限于：多模态/场景认知 agent、

心理咨询对话语料构建和过程建模等。

**科研目标：**

产出世界领先的学术成果，包括发表 AI 顶会论文 3 篇以上；举办或参加技术竞赛，扩大微信的技术影响力；定期举行技术讨论，产出成果在公司产品上进行落地应用，包括文本处理工具、翻译、搜索、对话等产品中的某些模块。

**可提供资源：**

GPU 计算集群，自然语言处理和对话数据集标注资源。

**关键词：**自然语言处理技术，人机对话技术，AI 认知技术。

## 5.2. 从文本到 SQL (Text2SQL) 的自动语义解析 (Semantic Parsing) 算法研究

微信支付的产品、运营、数据部门有频繁复杂的数据查询需求，如何高效支持这些需求，提升数据团队效能是非常重要的课题。Text2SQL 是一个值得探索的方向，Text2SQL 旨在研究如何将文字描述的数据需求自动翻译成机器语言 (SQL 查询语句)。2017 年来，不少 BI 工具已经提供了可以支持自然语言查询的功能，但都相对简单。聚焦于金融场景下，如何精确地支持更加复杂的查询，还需要在算法层面有更多突破。

本课题主要的研究内容主要聚焦于如何支持多表查询，以及如何准确匹配需求和库表字段。

**科研目标：**

研究自然语言到 SQL 语句的创新算法模型，提升在 Text2SQL 公共数据集上的效果，并将其落地为内部平台，提升效能。产出成果包括自然语言自动转化为 SQL 的机器学习模型，以及 1 篇顶会论文和 1 个专利。

**可提供资源：**

可提供部分脱敏数据和服务器资源等。

**关键词：**自然语言处理，机器学习。

## 6. 实验平台与软件测试

### 6.1. 实验平台算法模型研发

实验平台是互联网公司数据驱动的核心工具，用来敏捷迭代，快速试错，判断每一

个算法创新和产品创新是否对核心指标有显著正向提升，并且预测全量上线后的效果。本课题重点研究如下内容：

1、序列化检验研究：基于微信用户数据，做实时的序列化假设检验，随机过程的最优停止理论，Bayes factor，非IID数据下的optional stopping问题；

2、实验加速功能：研究加速实验的方法，在同等的检验精度下应用更少的样本，比如方差削减，优化指标，代理指标，强化学习调实验流量等方法的研究；

3、Bayesian optimization, Multi Armed Bandit在实验平台的优化，解决样本非IID，辛普森悖论等问题；

4、实验效果评估，如a. Causal Inference, heterogeneous treatment effects 传统的树模型假设强，效果不是很理想，结合实际数据特点，优化HTE相关的算法模型；  
b. Treatment effects估计，估计实验全量发部后的提升，比如经验贝叶斯方法估计。

**科研目标：**

运用统计学的方法，提升实验平台的准确性和时效性；产出论文，并将研究成果在生成环境中落地应用。

**可提供资源：**

测试数据集和计算资源。

**关键词：**统计，计算机。

## 6.2. 基于微服务依赖图的故障根因定位

许多公司已经将软件系统从单体架构迁移至微服务架构，并将大量核心业务实现为微服务系统。但微服务系统是一个十分复杂的系统，由于服务数量众多，服务之间的交互以及运行时环境的关系也变得极其复杂。这样的复杂性和动态性对微服务系统的故障定位带来了巨大而独特的挑战，这要求开发人员对整个系统的交互关系和拓扑结构有着全面的了解才能高效的定位系统故障并修复。但微服务系统通常会由不同团队开发，每个团队只负责自己模块的相关服务，难以建立对整个系统的全局了解，这都为定位故障问题带来困难。本课题计划研究自动化的微服务系统故障定位技术和工具，解决微服务系统故障定位困难的问题。

**科研目标：**

该项目计划研究基于微服务依赖图的自动化故障根因定位技术，基于服务调用链数据、服务日志数据、版本库数据等设计有效的自动化故障根因定位方法，使故障定位准确度和效率得到提升。项目预期实现原型工具一套，发表学术论文一篇，并在微信测试

团队以及相关团队推广使用。

**关键词：**软件测试， AIOps。