



YeeZ

Fidelius: 面向数据合作的隐私保护 区块链解决方案

Fidelius: YeeZ Privacy Protection for Data Collaboration —
A Blockchain based Solution

熠智科技

<https://yeez.tech>

2020 年 6 月

目录

1	背景介绍	1
1.1	数据合作	1
1.2	数据隐私保护	3
1.3	数据可用不可见	5
2	Fideliu s 解决方案	6
2.1	执行流程	7
2.2	系统架构	7
2.3	特性	9
3	Fideliu s 数据管理平台	9
4	Fideliu s 可信计算平台	11
4.1	智能合约	12
4.2	高性能数据分析库	13
4.3	静态隐私检查	13
4.4	EAnalyzer	15
5	Fideliu s 密态计算组件	15
附录 A	相关技术、产品	17
A.1	可信执行环境	17
A.2	定位于隐私保护的解决方案	18
A.3	结合了 TEE 的区块链系统	19
附录 B	为什么隐私保护需要区块链	20
B.1	什么是区块链	20
B.2	Fideliu s 中的区块链做了什么	20

1 背景介绍

1.1 数据合作

进入大数据时代，人类获取、管理和利用数据的能力空前提升，社会各界对数据的价值愈发重视¹。在数字经济时代，数据已经成为关键生产要素，就像在农业经济时代和工业经济时代中，土地、劳动力和资本是关键生产要素²。

相比传统生产要素，数据有着非常强的网络效应，当不同维度、不同来源的多种数据相结合的时候，其中蕴含的社会价值、经济价值往往能达到 $1+1>2$ 的效果。例如，医疗数据与位置信息的结合能够更快的控制传染病的发生规模；设备传感器的数据与用电数据的结合能够极大的提高用电效率；企业订单、物流数据与金融数据的结合能够为中小企业带来更精准的金融政策等。

可以说，在大数据时代，数据社会和商业价值的变现，需要多方参与，我们将这种为获得数据的综合见解而对多个独立数据进行的跨域数据处理称为数据合作³。

数据合作通常会涉及到多方参与，但一般可分为两类角色：数据提供方与数据使用方，前者提供数据，后者使用数据提供方的数据。

一般情况下在某个数据合作场景中可能存在多个数据提供方或者数据使用方，甚至某个主体即是数据提供方也是数据使用方。例如，一家在线旅游公司（OTA，Online Travel Agency）需要接入酒店预定的服务以优化旅游体验，作为数据使用方，该 OTA 公司会采购酒店预定服务提供商（数据提供方）的数据；而该 OTA 公司的用户数据也同时可以用于精准广告投放（此时该 OTA 公司则成为数据提供方）。

而更大范围的数据合作则通常需要引入中间人，后者帮助完成数据的展示、转发甚至必要的撮合、交易。中间人一般以中心化平台的形式呈现，近年来地方政府和产业界纷纷成立的大数据交易中心和平台⁴，正是承担着上述数据合作中间人的角色。图 1 描述了这两种数据合作的模型。

虽然企业、部门之间需要进行数据合作，但是在数据合作的实践上却差异很大，这是因为 1) 参与数据合作的主体之间的关系的的影响；2) 数据合作内容、目的的区别 3) 数据合作技术的限制等。

¹ 《开放的数林——政府数据开放的中国故事》，郑磊，上海人民出版社

² 《中共中央国务院关于构建更加完善的要素市场化配置体制机制的意见》，2019

³ What is data collaboration, <https://medium.com/infosum/what-is-data-collaboration-6519e604a365>

⁴ 全国大数据交易所及数据交易平台汇总, <http://www.tanmer.com/blog/541>

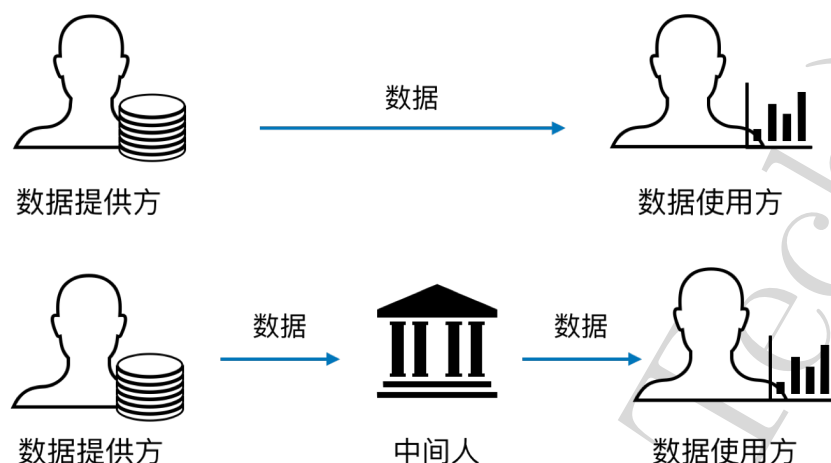


图 1: 数据合作模型

数据合作在表现形式上有多种形式，包括但不限于：

- 数据共享：通过某种方式实现多方原始数据的互通，例如彼此开放数据访问的 API、搭建专线、物理介质传输等；
- 数据交易：多方之间的数据交易行为，类似于数据共享，常见于存在中间人的情况，由中间人进行交易的撮合、数据的交付及支付；
- 数据分析：数据分析通常需要基于多方数据进行联合建模从而挖掘更多数据价值，例如近年逐渐受到关注的联邦学习。

数据合作的应用场景非常广泛，这里列举一些应用示例：

- 金融风控：银行在开展互联网贷款业务时，坏账率对于业务成败非常关键，因而需要对用户进行更精准的风控刻画，也需要针对整体的风控算法和参数，与外部多方协同合作；
- 营销分析：在广告营销领域，广告商和流量方需要接入第三方数据，从而实现精准定位客户来提升投放以及转化效果；
- 供应链金融：商业银行等金融机构通过结合供应链中的信息流、物流、资金流数据，来缓解金融机构和中小型企业之间信息的不对称问题，从而更有效的向核心企业以及供应商提供融资和贷款等业务；
- 园区综合能源服务：电网公司通过对园区内能源生产和消耗的数据的分析，可以实现包括发电和冷热负荷的协同利用，通过优化、经济分析，进而保证能源供应的安全性和经济性的平衡管理；

- 医疗数据共享：医疗信息对保险公司、医保机构和流行病监测部门具有至关重要的作用，医疗信息的共享可极大地提高保险的精准化水平、医保部门政策的精细化和流行病防治的及时性；同时诊疗流程中有极多的信息点可以帮助后期的临床质量提升或新治疗方案的产生，医疗机构可以通过共享医疗数据提供临床辅助决策支持、应用管理决策支持等；
- 共享征信：当前我国获取主体信用信息的模式有数据中心模式、第三方征信模式、共享查询模式三类。传统征信机构通常采用数据中心模式，出于对个人信息的保护，央行往往审慎从严下发个人征信牌照，类似地第三方征信机构利用自身系统或技术优势，构建中心化平台向授信机构提供服务；而近年来的共享征信则是数据合作的典型应用：业务机构无需事先将数据上报给共享中心，数据由机构自行管理，当机构需要获取数据时，通过中心发送到其他机构，有数据的机构回应信息，返回查询机构。

当前业界已经逐渐意识到数据合作是数据作为生产资料的必然诉求⁵，然而，如何保证数据提供方的数据产权及如何保护数据的隐私是企业、部门之间开展数据的关键前提。

1.2 数据隐私保护

隐私保护是一个十分庞杂的概念，通常是指使个人或集体等实体不愿意被外人知道的信息得到应有的保护。对于个人来说，一类重要的隐私是个人的身份信息，即利用该信息可以直接或者间接地通过连接查询追溯到某个人；对于集体来说，隐私一般是指代表一个团体各种行为的敏感信息⁶。不难看出，对于个人和集体（企业）而言，隐私保护的概念和范畴不尽相同。

对于个人而言，隐私数据是自身和周边相关环境的个人数据，同时个人隐私数据的使用受到严格监管，这不仅仅依靠信息安全技术来解决，更多的是依赖于相关法律法规的约束。

欧盟《通用数据保护条例》（General Data Protection Regulation，简称 GDPR）定义了任何收集、传输、保留或处理涉及到欧盟所有成员国内的个人信息的机构组织所需要受到的约束；而刚刚于 2020 年新颁布的《中华人民共和国民法典》，以“隐私权和个人信息保护”专章方式，对隐私权和个人信息定义、保护原则、法律责任、主体权利、信息处理等问题作出规定。

⁵ 《蚂蚁金服在大数据合作上的创新实践》，ArchSummit 2016

⁶ <https://wiki.mbalib.com/wiki/隐私保护>

对于企业而言，情况则更为复杂。尽管企业数据包含了其收集的个人用户的数据（例如移动 APP 用户信息、电商平台用户数据、保险公司用户数据），但企业数据也包括了自身和合作伙伴的关键商业数据（例如经营收入、经营利润、市场占用率等）以及业务数据（例如产品核心数据、生产线数据等）。上述数据既属于企业隐私，也是企业的重要资产，因此，这类数据的泄露会进一步造成企业隐私的泄露；同时，由于数据具有容易复制、不易追踪的特点，其边际使用成本极低，而可能出现二次贩卖会减少企业在数据合作中的收益，也会削弱企业核心竞争力。

隐私泄露问题越来越引起社会重视，对于企业而言，内部员工利用职责便利泄露数据等行为一直难以避免⁷。然而，即便看似“安全”的第三方数据合作平台，在利益驱使下也可能出现数据倒卖的问题⁸。

因此，数据隐私保护是企业之间、部门之间开展数据合作的关键前提，除了监管和相关法律法规的制定等制度层面的措施之外，也应该加强技术手段从根本上杜绝隐私泄露的可能。

值得注意的是，隐私保护也并非单一的技术方案，从数据的发布到数据存储、数据使用过程中，涉及到不同类型的隐私保护技术。

目前主要的隐私保护技术通常包括：

- 数据发布隐私保护：在数据发布时，需要保证用户数据可用的情况下，高效、可靠地去掉可能泄露用户隐私的内容。k-匿名 [?]、l-diversity 匿名 [?] 等针对数据的匿名发布技术可以实现发布数据时的匿名保护。
- 数据存储隐私保护：在数据存储在云平台或者第三方存储服务商时，存储并不能保证是完全可信的。用户的数据面临着被不可信的第三方偷窃数据或者篡改数据的风险。对称加密技术 DES、AES [?] 以及不对称加密技术 RSA [?]、Elgamal [?] 等是解决该问题的传统思路，此外混合加密技术、同态加密技术 [?] 等也是针对数据存储时防止隐私泄露而采取的一些方法。
- 数据使用隐私保护：在数据的使用（例如数据查询、数据分析等）过程中，需要在尽可能保证数据可用性的前提下，采用合适的数据隐藏技术，以防范利用数据发掘方法引发的隐私泄露。现在的主要技术包括：基于数据失真和加密的方法，比如差分隐私 [?] 以及多方安全计算（MPC） [?]、可信执行环境（TEE） [?] 等技术。

⁷ 《实名举报！脱口秀演员池子起诉中信银行：未经授权泄露个人隐私》，新浪财经

⁸ 2018 年大数据交易平台“数据堂”被查出倒卖 4000GB 的隐私数据给 Google、三星、Microsoft 等境外企业牟取暴利

1.3 数据可用不可见

在企业的数据合作场景中，由于数据分布在多个主体中，并且数据的使用权和拥有权通常并非同一主体，因此数据的隐私保护也更为复杂。

一方面传统的隐私保护技术存在缺陷无法完全避免隐私泄露⁹；另一方面，数据提供方将数据（即便经过匿名处理）交付给数据使用方后，则丧失了数据的拥有权，例如数据的二次贩卖等行为将难以阻止。

实际上对于很多数据使用方而言，其本身诉求并非原始数据，而是基于数据分析做出数据驱动决策（data-driven decision-making）[?]。例如在供应链金融场景中，金融机构并不关心供应商实际的物流、销售数据而是需要基于上述数据判断后者是否符合融资或者贷款要求。这为数据合作提供了新的解决思路。

一个直观的思想是，通过将对于数据的使用或者计算进行迁移，数据提供方提供数据处理服务而非原始数据，可以从根本上避免数据出域造成隐私泄露，上述思想通常被称为“数据可用不可见”。

目前，业界实现数据“可用不可见”的技术路线主要有两条：

1. 基于密码学技术的密态计算：以安全多方计算、可搜索加密、同态加密、零知识证明等技术为代表。其核心思想有二：其一，数据提供方本地使用原始数据计算分析结果的同时生成一个基于零知识证明技术的证明。该证明只有在分析结果正确的前提下才能够被生成，同时不会泄露原始数据的任何信息。这样能让数据不出域的同时保证正确性；其二，设计特殊的加密算法和协议，从而支持在加密数据之上（不用解密）直接进行计算，得到所需的计算结果，同时不接触数据明文内容。
2. 基于可信执行环境技术（TEE, Trusted Execution Environment）的可信计算：以 Intel 的 SGX, AMD 的 SEV, ARM 的 Trust Zone 等技术作为代表。其核心思想是以可信硬件为载体，提供硬件级强安全隔离和通用计算环境，在完善的密码服务加持下形成“密室”，数据仅在“密室”内才进行解密并计算，除此之外任何其他方法都无法接触到数据明文内容，数据在离开“密室”之前又会被自动加密，从而实现“可用不可见”。

“数据可用不可见”的思想理论上可以有效实现数据合作中的隐私保护，目前也有一些分别基于上述两种机制的隐私保护解决方案¹⁰，但是我们发现仍然存在一些关键问题阻碍了这一技术的实践：

⁹例如经过匿名等处理后的数据，通过大数据关联分析、聚类、分类等数据挖掘方法，依然可以分析出部分隐私数据

¹⁰关于相关技术和产品介绍参见A

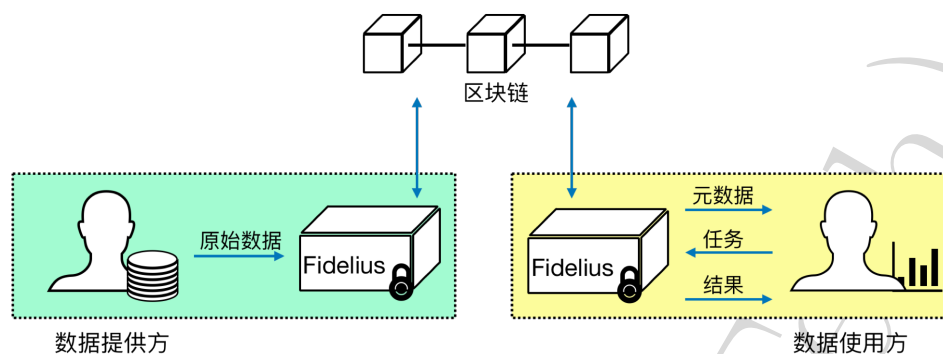


图 2: 隐私保护数据合作模型

1. 原始数据的一致性：即如何保证计算过程的输入数据是数据提供方声称提供的数据，且没有被篡改；
2. 计算逻辑可控：即如何保证结果结果不包含数据提供方意愿之外的信息，从而造成数据不可见情况下的隐私泄露；
3. 计算结果的正确性：即如何保证计算结果是由数据使用方提供的计算程序（且未经篡改）生成的，而不是数据提供方随意生成的；
4. 计算结果的隐私性：即如何保证计算结果仅数据使用方可见，数据提供方不能查看计算结果。

这些问题在数据使用方获得原始数据的情况下都是天然保证的，而在“数据可用不可见”的情况下，这些问题的重要性便凸显出来，甚至在一定程度上阻碍了数据合作的可能性。

2 Fidelius 解决方案

为了赋能企业间的数据合作，助力企业利用数据提升自身核心竞争力，熠智科技推出了面向数据合作的一站式隐私保护解决方案 Fidelius。Fidelius 基于“数据可用不可见”思想，同时有效的保证了原始数据的一致性、计算逻辑的可控性、计算结果的正确性及隐私性。

图 2描述了基于 Fidelius 实现数据合作的抽象流程。与传统的数据合作模式类似，参与方包括了数据提供方和数据使用方。为了简明，我们忽略数据真实的存储位置，在某些系统设定中，将其统称为在线数据仓库¹¹。

¹¹Solid, <http://solid.mit.edu>

FideliuS 中间件分别运行在数据提供方和数据使用方中，双方通过与 FideliuS 交互实现数据合作操作。数据提供方和数据使用方之间没有直接的数据交互，并且原始数据不会离开数据提供方的 FideliuS 中间件，这从根本上避免了隐私数据的泄露问题。

值得注意的是，相比图 1 所示的传统的数据合作模式，FideliuS 引入了区块链网络。由于区块链本身具有去中心化网络、公开可验证等特性，FideliuS 将其作为可信的传输通道和数据计算验证平台¹²。

2.1 执行流程

尽管数据合作表现形式有多种，但仍可以将其核心流程抽象为图 2，在此基础上我们简述 FideliuS 的执行流程：

1. 数据注册：数据提供方传输原始数据至本地的 FideliuS 组件（此处简称 FP），后者自动生成相关元数据（即原始数据的描述信息）和权限证明，同时元数据被发布至区块链网络；
2. 任务发布：数据使用方所在的 FideliuS 组件（此处简称 FC）通过区块链获取元数据后提交给数据使用方，后者根据数据描述信息提供相应的数据分析任务（通常为执行程序的二进制文件），FC 将数据分析任务通过区块链网络转移到至 FP；
3. 数据计算：FP 对任务进行检测后，基于原始数据和任务完成数据计算，生成计算结果和相应的证明，并将加密后的计算结果和证明发布至区块链网络；
4. 结果返回：分析结果（加密）被保存在区块链网络中，区块链中的智能合约基于证明和结果进行验证，当计算结果通过验证后，FC 将计算结果解密后返回给数据使用方。

2.2 系统架构

FideliuS 包括了数据管理平台、隐私计算平台以及密态计算组件，各个部分包含的模块以及关系如图 3 所示。

数据管理平台 数据管理平台为隐私计算平台及密态计算组件提供数据源，功能上包括了隐私数据的注册和权限管理，数据的注册模块负责元数据生成和加密；权限管理

¹²关于区块链在 FideliuS 中功能的介绍，参见附录 B

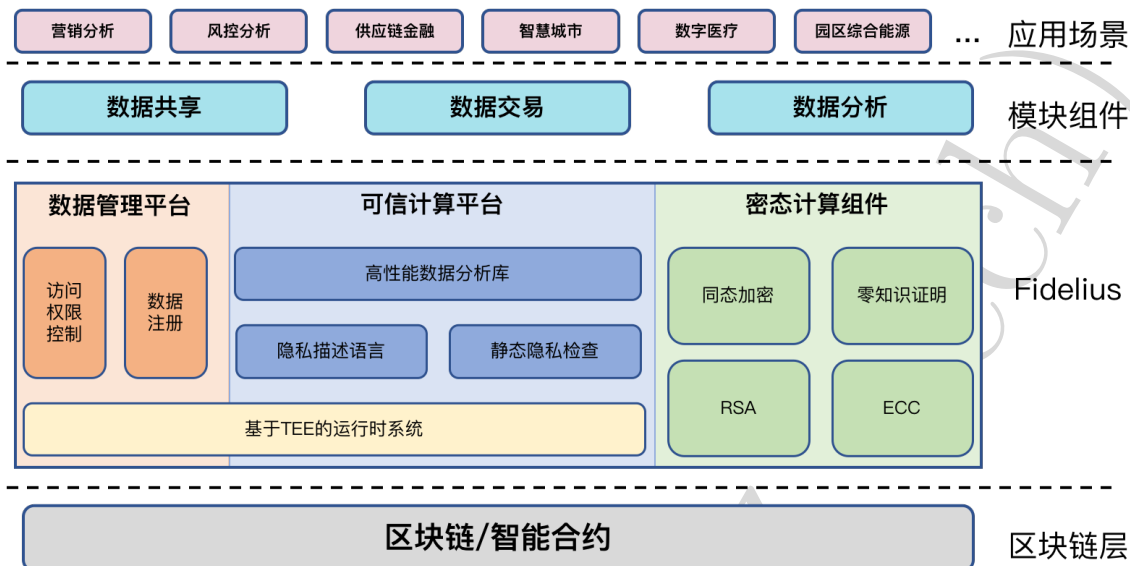


图 3: FideliuS 架构

模块定义了原始数据的访问方式、格式、权限等信息。需要注意的是，数据管理平台不影响原始数据的存储方式或介质，仅在数据使用过程中仅临时访问原始数据并进行加密。

隐私计算平台 负责完成数据使用的交互过程，包括数据分析任务的生成、发布、执行，以及计算结果的验证以及交付。隐私计算平台基于 TEE 可信执行环境，提供了数据的通用计算分析支持。特别的，隐私计算平台确保了原始数据的一致性、计算逻辑的可控性、计算结果的正确性及隐私性。

密态计算组件 主要负责数据分析任务的加速。对于分析任务中计算较为繁琐的部分，密态计算组件支持将该部分计算脱离于 SGX 之外（如在 CPU，GPU 上）进行计算，如支持并行运算等，从而具有更高的执行效率。同时，密态计算组件仍然确保原始数据和中间结果的隐私性、执行过程的正确性。

此外 FideliuS 底层实现了对区块链的交互接口，后者负责元数据的存储、可信传输和公开验证。FideliuS 并不依赖于特定的区块链实现，因此能够部署在不同区块链系统上。当然，FideliuS 的性能，如吞吐量、响应时间等，受到区块链系统的影响，我们推荐使用高性能的 YeeZChain 作为底层区块链平台¹³。

¹³关于 YeeZChain 详细介绍参见《YeeZchain v1.1 介绍文档》，https://gitlab.com/Yeez/yeetz_introduction

2.3 特性

相比于传统数据合作方式，FideliuS 消除了数据提供方对于隐私泄露的担忧：对于数据提供方而言，由于数据使用方无法接触原始数据，同时计算逻辑可控制，因此无需担心数据隐私泄露。

同时对于数据使用者来说，原始数据的真实性和任务执行的正确性可以得到验证，因此真正做到了“数据可用可信不可见”。

更具体来说，FideliuS 赋能下的数据合作主要有以下特点：

1. 覆盖数据全生命周期的隐私保护：即“不可见”，FideliuS 实现了整个数据合作流程中数据的隐私保护：首先，FideliuS 不直接存储原始数据；在数据注册过程中，FideliuS 仅发布了原始数据的描述信息，即元数据；而在数据计算过程中，FideliuS 通过静态检查进一步确保了数据分析任务不会违背数据提供方的意愿泄露隐私数据。
2. 高性能通用隐私计算：即“可用”，FideliuS 将数据计算与数据存储进行了解耦，由隐私计算平台和密态计算组件负责数据的计算流程。隐私计算平台基于 TEE 可信执行环境提供了高性能数据分析库，支持包括聚类、分类、回归以及关联分析等常用的数据分析工具。而对于不依赖于 TEE 的计算需求，则可以基于密态计算组件实现，后者可以基于通用的 CPU 或 GPU 并行加速完成。不管是隐私计算平台，还是密态计算组件，都能够应对大规模的数据处理。
3. 数据分析结果可验证：即“可信”，这主要得益于两个方面：1) 数据提供方不能篡改数据或提供不完整的数据；2) 数据提供方不能篡改数据使用方的数据分析任务；同时区块链之上的智能合约会在计算结果的验证阶段对上述信息进行公开验证，确保交付的数据分析结果的可信性。

3 FideliuS 数据管理平台

FideliuS 数据管理平台（简称 F-DMP）适用于数据访问需要进行严格管理的场景，在这些场景下，F-DMP 能够有效的消除数据访问权限管理不当带来的系统性风险。相比较而言，传统的数据访问权限管理是存在系统性风险的，即存在超级管理员，超级管理员能够绕开所有的权限接触数据或者在不留下痕迹（删除特定的访问记录）的情况下接触数据。2018 年 3 月曝出的 Facebook 超过 5000 万用户信息数据遭到外泄¹⁴，2016 年 11 月爆出的“陕西千亿矿权案”的案件卷宗在最高人民法院办公室

¹⁴http://www.xinhuanet.com/world/2018-03/24/c_129836684.htm

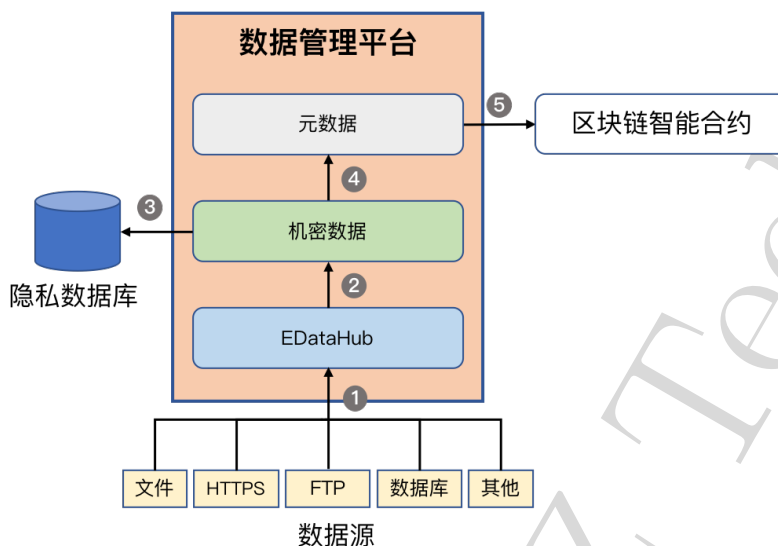


图 4: F-DMP 发布数据工作流程

失窃¹⁵都表明了传统权限管理的系统性风险。

F-DMP 通过 TEE 和区块链技术的结合，将数据的元数据发布在区块链智能合约上。图 4描述了在 F-DMP 中发布数据的工作流，其中 EDataHub 运行在 TEE 中，保证了对数据访问的可信性；机密数据包括了访问数据源的机密信息，包括用户名、密码、路径、证书等信息，这些信息与数据所在的环境有关，不能发布到区块链上，因此仅能存储在本地，更进一步的，这些机密数据使用了 EDataHub 特有的加密密钥进行加密，以保证这些数据的机密性；元数据则包括了数据的名称、条数、格式、权限要求、哈希等信息；最后，区块链智能合约中记录了相应的数据元数据。F-DMP 能够通过智能合约定义丰富的权限管理语义，例如，基于密级的权限控制、基于访问记录的动态权限控制、基于角色的权限控制、基于属性的权限控制等。

在 F-DMP 中，对数据的访问或使用需要通过智能合约发起，如图 5所示。在访问或使用数据的过程中，同样需要使用 EDataHub 访问相应的数据。在使用 EDataHub 得到原始数据后，如果是需要根据访问权限进行数据访问，则需要对原始数据进行加密、签名后提交到智能合约；如果是需要对原始数据进行数据分析，则对原始数据进行加密后，转交给本地的、运行在 TEE 中的分析程序进行处理。

F-DMP 本身并不进行数据的存储，只是对元数据、访问权限及访问记录的管理。对于数据的存储，F-DMP 依赖于已有的数据源，包括文件、HTTP/HTTPS、FTP、各种数据库以及其他的数据存储系统。对于某些数据存储系统，如 HTTPS，已经包括了数据加密访问的机制，而对于其他不包括数据加密访问机制数据源、如文件、FTP 等，则需要根据需要进行进一步的处理。

¹⁵<https://zh.wikipedia.org/wiki/陕西千亿矿权案>

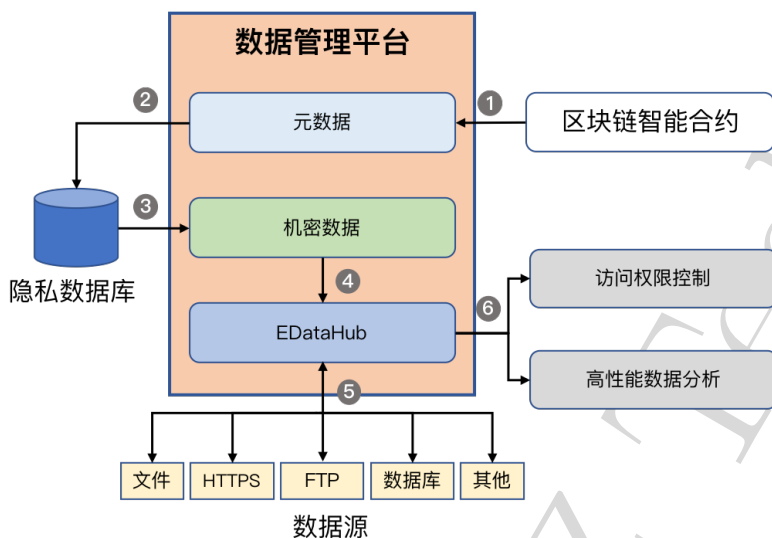


图 5: F-DMP 数据管理平台使用数据 workflow

综上，F-DMP 从三个方面消除了传统数据访问权限管理中的风险：

1. 区块链记录所有的访问记录，保证了数据访问记录的不可篡改；
2. 原始数据的访问过程（加/解密）在 TEE(Intel SGX) 中进行，保证了即使能够接触到物理机，也无法直接访问数据；
3. 数据在传输过程中使用访问者的公钥加密，保证了数据访问的隐私性。

4 Fidelius 可信计算平台

在 Fidelius 可信计算平台（简称 F-TCP）上，用户分为两个角色，Data Provider 和 Data Analyzer，即前述的数据提供方与数据使用方。Data Provider 通过 F-DMP 发布指定的数据，并响应 Data Analyzer 的计算、分析请求；Data Analyzer 则从区块链上查看元数据，并发布相应的计算、分析请求。

如图 6 所示，Data Analyzer 首先从区块链上查看元数据，获取数据的相关信息；然后根据这些元数据编写数据分析任务，在这一过程中，F-TCP 提供了相应的编程库及开发工具，在完成数据分析任务的编写及调试后，F-TCP 会进行静态隐私检查，以确保 Data Analyzer 的数据分析不会泄露 Data Provider 在意的隐私；在完成静态隐私检查后，F-TCP 生成 EAnalyzer，一个运行在 TEE 环境下的二进制数据分析程序。EAnalyzer 作为数据分析请求的一部分被 Data Analyzer 提交到区块链智能合约上。

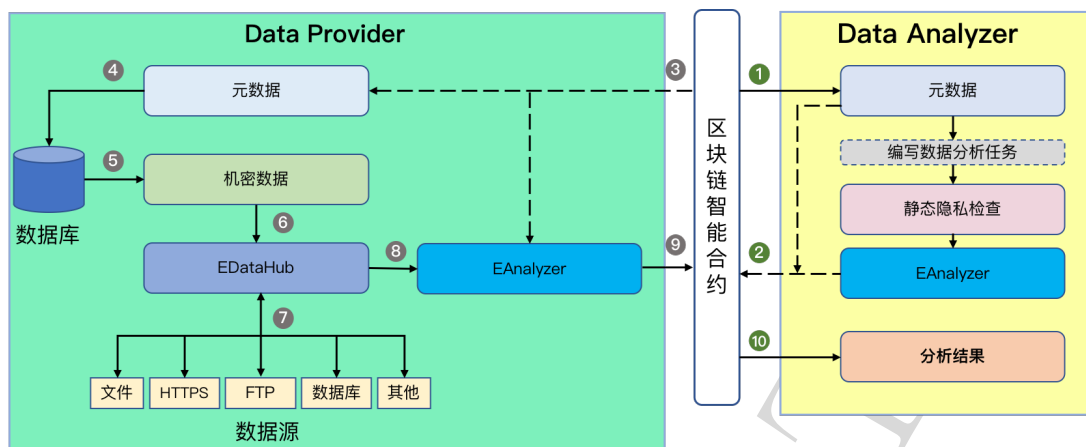


图 6: YeeZ 隐私计算平台 workflow

Data Provider 通过区块链获取数据分析请求之后，通过 F-DMP 中所述的流程读取原始数据，并将原始数据转交给本地运行的 EAnalyzer，EAnalyzer 执行完成，将分析结果加密、签名后上链。最后，Data Analyzer 从区块链上获取数据分析结果，并使用自己的私钥解密分析结果。

如图 6 所示，整个 workflow 涉及到多个模块，包括智能合约、高性能数据分析库、静态隐私检查、以及 EAnalyzer，下面分别介绍这些组件。

4.1 智能合约

对于每一个发布的数据项，Data Provider 需要在区块链上创建一个新的智能合约实例，我们使用区块链上的智能合约完成三个功能

- 数据元信息的发布：数据项元信息会被记录在合约中，如数据的 hash，名称，描述，条数，格式，价格以及其他必要的信息，关键信息之外的信息，如名称、描述，可以在发布后进行修改，修改权限可以通过智能合约灵活指定；
- 分析程序的记录、传输及请求：Data Analyzer 将自己的分析程序及其他信息发送给合约，并向合约支付指定的代币，智能合约记录此次分析请求；
- 分析结果的验证：任何人都可以向智能合约提供一个分析结果，智能合约会验证分析结果的合法性，仅当分析结果合法时，智能合约接收结果，完成数据分析任务。

4.2 高性能数据分析库

高性能数据分析库是一个运行在 TEE 环境下的数据分析库，使用 C++ 开发。高性能数据分析库使用基于数据流的计算模型，不同的计算单元按照计算任务的需求组成有数据流图。F-TCP 内置的计算单元的包括了对数据的过滤、拆分、降维、数据增强等基本操作，也包括了对数据集合的处理，例如聚类、分类、回归、关联分析等复杂的算法。

在数据流模型之外，高性能数据分析库还支持联邦学习。联邦学习的核心思想即在保证组织间信息交换过程中用户隐私和数据安全的前提下，实现各数据源在逻辑概念上的整合，从而构建全局模型，其被广泛应用于金融、医疗和教育等行业。然而，由于传统的联邦学习实现时不可避免的会存在一个逻辑的中心，从而使得其依然会存在信息泄露的风险。而基于 F-TCP 实现的联邦学习框架，一方面可以依托区块链实现模型的可信传输通道；另一方面还可以实现模型训练的“按劳分配”，即，依托 F-TCP，可以对各个数据提供方的贡献值进行计算，从而使得联邦学习训练过程中的全局记账和激励分配的可验证。此外，传统的联邦学习方法还极易受到后门攻击，而 F-TCP 提供的不可篡改和可追溯性有助于实现篡改和恶意行为的及时检测和替换，从而增强联邦学习模型的可审计性。最后，F-TCP 上模型的参数更新都是通过加密方式进行的，也进一步提高了模型的完整性和保密性。

受限于 TEE 的执行环境的影响，高性能分析库在内存分配、计算任务调度方面进行了针对性的优化。不同于适用于 CPU、GPU 的数据分析库，TEE 执行环境有着内存受限、上下文切换代价高的特点，因此，F-TCP 通过计算任务的调度尽量增加了内存的复用，减少数据的移动及计算任务的上下文切换，从而使得高性能分析库能用于大规模的数据分析。

4.3 静态隐私检查

静态隐私检查是 F-TCP 通过静态程序分析的手段，确保 Data Analyzer 编写的数据分析任务不会泄露 Data Provider 的隐私。一方面，Data Analyzer 的分析程序是不受 Data Provider 控制的，而这些不受控制的数据分析是有可能反映某些隐私相关的信息的，更有甚者，数据分析的结果可能直接使用原始数据做为数据分析的结果；另一方面，很难广泛的、通用的定义如何使用一个数据才构成了隐私泄露。

例如，对于用户名、密码之类严格的数据，任何信息的泄露都可能构成隐私泄露：即使是密码长度这一信息的泄露，也对信息安全造成了巨大的威胁；相对的，一个温度传感器在当天采集到的温度的原始数据可能是敏感的（这可以反映周围的设备的工作、负载状况），但平均温度或峰值温度是否超过 45°C ，则是可以被 Data Analyzer 获取的。

```
import "user_type.h"
import number;
in = input user_type_t;
iris = in.iris_data;

sum(s1, s2) = {data: s1.data + s2.data, counter: s1.counter + s2.counter}

s1 = {data: iris.sepal_len, counter:number.one} |
    sum(s1:s1, s1:s2);

sw = {data:iris.sepal_wid, counter:number.one} |
    sum(sw:s1, sw:s2) ;

pl = {data:iris.petal_len, counter:number.one} |
    sum(pl:s1, pl:s2) ;

pw = {data:iris.petal_wid, counter:number.one} |
    sum(pw:s1, pw:s2) ;

osl = s1.data/s1.counter;
osw = sw.data/sw.counter;
opl = pl.data/pl.counter;
opw = pw.data/pw.counter;

output osl, osw, opl, opw, in.species
```

图 7: 使用 YPDL 描述允许在输入数据上执行 KMeans

上述例子说明了隐私保护的规则需要 Data Provider 指定，也就是说，Data Provider 需要描述数据应该被如何使用，为此，我们引入了隐私描述语言（PDL）。PDL 将数据的运算规则描述为有限状态机，反映了从输入数据到输出数据的状态转换，因此，不在 PDL 描述的状态转换（运算），都会被认为违反了 Data Provider 定义的隐私规则。图 7 描述了一个使用 PDL 描述一个数据可以进行 KMeans 聚类分析的示例。

F-TCP 使用了静态二进制分析的方法检查用户的二进制程序是否遵守 PDL 的描述。在静态二进制分析中，使用了 GTIRB（The GrammaTech Intermediate Representation for Binaries）¹⁶作为分析的中间表示，在将二进制程序转换为中间表示后，PSC 使用符号执行（Symbolic Execution）获取每个输出变量的状态，并将这一状态与 PDL 的描述进行匹配，如果不匹配，则认为存在隐私泄露。

为了保证静态隐私检查结果的可靠性，即避免 Data Analyzer 篡改隐私检查的过程或结果，静态隐私检查运行在 TEE 上，当且仅当静态隐私检查通过时，才会生成最终的 EAnalyzer。

¹⁶<https://github.com/GrammaTech/gtirb>

4.4 EAnalyzer

EAnalyzer 包含了 Data Analyzer 指定的数据分析任务，运行在 TEE 之上，以保证计算过程不可见、不可篡改。EAnalyzer 基于指定数据的元数据进行编写，主要包括数据的格式，并且 EAnalyzer 对数据的使用需要满足 Data Provider 提供的隐私描述所提供的规则；EAnalyzer 的开发基于高性能数据分析库，可以使用数据流的方式对数据进行分析，也可以实现联邦学习任务。

EAnalyzer 运行在 Data Provider 的环境下，Data Analyzer 仅仅编写、生成了 EAnalyzer，因此 EAnalyzer 是在数据所在的域运行，保证了数据不出域的特性。EAnalyzer 的数据来源是通过 EDataHub 得到，EAnalyzer 与 EDataHub 之间使用 ECDH 交互的密钥加密交换的数据。

EAnalyzer 首先需要保证计算结果的正确性与机密性。对于正确性，是指 Data Provider 提供的数据分析结果严格执行了 EAnalyzer 得到的，而不是随意生成的；对于机密性，则是指 EAnalyzer 的计算结果不能被相应的 Data Analyzer 之外的用户获悉。为了达到这两个目标，我们将只有 Data Analyzer 知道的私钥加密后随 EAnalyzer 发送到 Data Provider，且仅有该相应的 EAnalyzer 能使用该私钥。在 EAnalyzer 完成计算分析后，使用该私钥对计算结果进行加密、签名。由此，区块链智能合约能够通过该签名验证数据是否是由指定的 EAnalyzer 生成，从而保证了计算结果的正确性；而计算结果的机密性则是由于智能合约接受的数据为密文，只有相应的 Data Analyzer 能够使用自己的私钥解密相应的计算结果，如图 6 所示。

EAnalyzer 还需要保证输入数据的一致性，即计算中所使用的数据与 Data Provider 所声称的数据是一致的，既没有被篡改也没有遗漏。EAnalyzer 在计算分析的过程中，对每个输入数据的条目都进行一次哈希运算，随着计算过程的结束，得到最终的输入数据的哈希。EAnalyzer 使用前述的私钥（只有相应的 Data Analyzer 知道，并随 EAnalyzer 发送到 Data Provider）对这一哈希进行签名，并由区块链智能合约对该签名进行验证，从而保证了输入数据的一致性。《《《《 Updated upstream

5 Fidelius 密态计算组件

Fidelius 密态计算组件（简称 F-CCM）从另外一个（纯软件）角度解决了数据合作问题，不需要任何一方的 TEE 组件作为支撑。如图 8 所示，DA（Data Analyzer）可以直接将静态分析程序发布至区块链上（该分析程序不需要 SGX 做封装。当然也需要通过 YPC 的静态隐私检查），同时发布相关数据分析请求。DP（Data Provider）一旦接收到数据分析请求以及相应的静态分析程序之后，在本地调用原始数据并执行分析程序，得到程序执行结果。同时，DP 需要（同样本地）生成一个（零知识）证

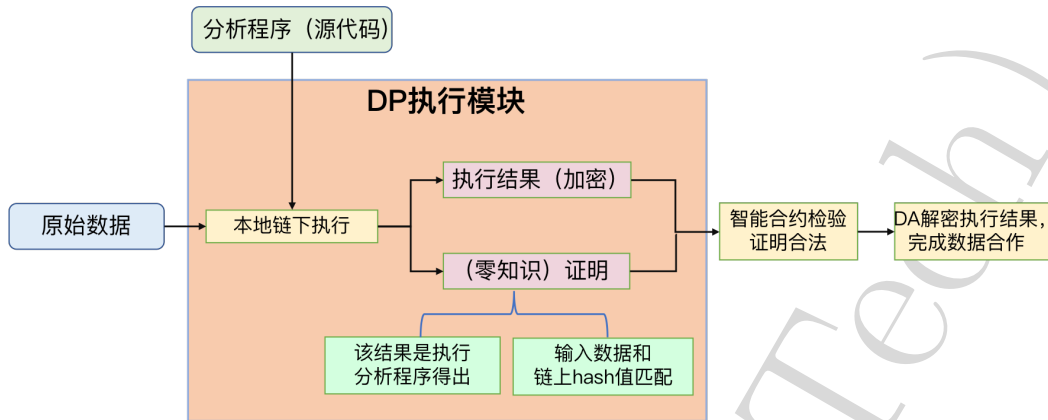


图 8: 解决方案示意图

明。该证明基于密码学技术，能够保证 1) 该（用 DA 公钥加密之后的）执行结果确实是经由 DA 提供的分析程序得到。2) 分析程序的输入（对应 DP 的原始数据）的 hash 值和记录在合约中的原信息匹配。随后，DP 将程序执行结果用 DA 公钥加密后，附带生成的证明一起发送到链上，交由智能合约验证。智能合约验证 DP 提供的证明满足上述两点保证，该过程完全公开。验证通过后，DA 即用自己的私钥对执行结果进行解密，得到想要的分析结果，此次数据合作完成。

F-CCM 保证数据合作安全可用的关键在于生成的这个证明，其具有如下性质：

- 可靠性：如果 DP 正确的执行了分析程序且输入数据合法（原始数据的 hash 值与链上匹配），则生成的证明能够通过。反正，如果 DP 没有按照给定的分析程序来执行，或者 DP 没有给出合法的输入数据，则 DP 难以伪造出一个合法的证明来通过验证。
- 非交互性：DP 可以本地直接生成合法证明，其过程中不需要和 DA 进行额外交互，
- 零知识性：该证明不会泄露 DP 的任何隐私信息，包括原始数据信息，中间变量信息等。
- 简洁性：该证明的长度通常为常数级别，且验证其合法性的时间复杂度远远小于重新执行分析程序，能支持智能合约验证。
- 通用性：支持对绝大部分高级语言分析程序生成证明。

在 Fidelius 的设计中，TEE 和 F-CCM 可以结合起来使用实现效率最大化：通常，分析程序中逻辑复杂的部分（如大量条件分支语句，大量递归嵌套等）不适合使用密码学组件，因为此部分代码生成的证明往往较为复杂。故此部分我们推荐用 TEE 做隐私保护。而逻辑简单，但运算繁琐的部分（能加速的操作，数据并行操作，

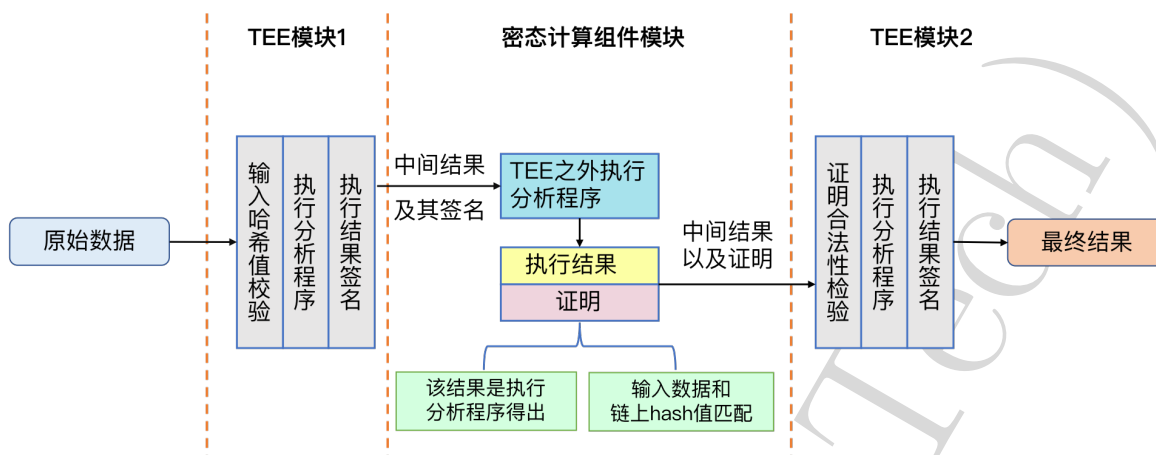


图 9: 解决方案示意图

大规模矩阵乘法运算）我们可以置于 TEE 之外完成，因为在 CPU 甚至 GPU 上执行运算相比于 TEE 更有效率，且生成的证明也极其简洁。

其具体的系统架构如图9所示，其中，每个模块的输出当做下一个模块的输入。同时，对于 SGX 模块，除第一个模块外，对输入哈希值的验证改为验证输入及附带证明的合法性（由上一个密码学组件模块的输出提供）。而对于密码学组件模块，其生成的证明应包括对输入及附带签名的验证（由上一个 SGX 模块的输出提供）。此架构允许在 SGX 执行程序中插入任意数量的密码学组件以实现加速，同时仍然只需和区块链进行一次交互即可保证正确性。

附录 A 相关技术、产品

整个隐私保护产业目前仍处于早期发展阶段。本章简单介绍下业界已有的相关技术和解决方案。

A.1 可信执行环境

可信执行环境（TEE, Trusted Execution Environment），它是电子设备（PC、智能手机、平板电脑等设备）CPU 上的一块区域。这块区域的作用是给数据和代码的执行提供一个更安全的空间，并保证它们的机密性和完整性。TEE 最早来自于移动终端开放组织（OMTP）在 2006 年提出一种解决方案：在同一个智能终端下，除了多媒体操作系统外再提供一个隔离的安全操作系统，这一运行在隔离的硬件之上的隔离安全操作系统，用来专门处理敏感信息以保证信息的安全。

目前提供 TEE 产品的主要来自英特尔、ARM 等芯片厂商。

- SGX: Intel SGX¹⁷是 Intel 架构新的扩展, 在原有架构上增加了一组新的指令集和内存访问机制。这些扩展允许应用程序实现一个被称为 enclave 的容器, 在应用程序的地址空间中划分出一块被保护的区域, 为容器内的代码和数据提供机密性和完整性的保护, 免受拥有特殊权限的恶意软件的破坏。
- TrustZone: TrustZone¹⁸是 ARM 公司提出的一种 TEE 实现方案, 其将 SoC 的硬件和软件资源划分为安全和非安全两个区域。当设备执行隐私相关的操作时, 比如指纹识别、密码处理、数据加解密、安全认证等, 会在安全区域执行, 其余操作在非安全区域执行, 比如用户操作系统、各种应用程序等。

A.2 定位于隐私保护的解决方案

基于 TEE 技术, 许多厂商提出了隐私保护产品和解决方案, 代表性的企业 and 产品主要包括:

- Fortanix: Fortanix¹⁹是使用 Intel SGX 的 TEE 技术较早的公司之一, 其产品主要包括自防护密钥管理服务 SDKMS, 用于保证密钥的生成和使用的安全性; 以及用于密文计算的运行时加密平台 Runtime Encryption, 其允许用户自定义加密逻辑并在英特尔 SGX 的可信执行环境中运行。
- Asylo: Asylo²⁰是谷歌基于 Intel SGX 技术推出的开源安全应用开发框架, Asylo 框架可以使开发人员能够方便地使用 TEE 技术实现密文计算, 并支持从企业内部系统到云端的部署。
- MesaTEE: MesaTEE²¹是百度基于 Rust 语言和 Intel SGX 技术推出的通用安全计算框架, 其宣称采用了混合内存安全技术、密文计算技术, 以及可信计算技术 (如 TPM), 通过私有化或云服务帮助金融、政务、互联网等行业在联合建模、联合营销、联合风控等场景下一站式完成数据联合计算。
- 蚂蚁摩斯: 蚂蚁摩斯 (ANT MORSE)²²是蚂蚁金服结合区块链和安全计算技术推出的数据安全共享基础设施, 其目标为解决企业之间数据合作过程中的数据安全和隐私保护问题, 打通数据孤岛。从产品定位而言, 蚂蚁莫斯和 Fidelius 较为一致, 但在技术实现上, Fidelius 相比于蚂蚁摩斯能够保证计算结果的正确性和机密性。

¹⁷<https://software.intel.com/content/www/us/en/develop/topics/software-guard-extensions.html>

¹⁸<https://developer.arm.com/ip-products/security-ip/trustzone>

¹⁹<https://fortanix.com/>

²⁰<https://asylo.dev/>

²¹<https://anquan.baidu.com/product/mesatee>

²²<https://tech.antfin.com/products/MORSE>

此外，也有基于密码算法、多方安全计算实现的隐私保护解决方案，例如：

- zk-SNARK: zk-SNARK²³是 zero-knowledge succinct non-interactive arguments of knowledge 的简称，zkSNARK 号称使用了简洁化的非交互式零知识证明，旨在实现通用的隐私保护协议，目前 zk-SNARK 已被广泛应用于以太坊等数字加密货币中，类似实现的还有 zk-STARK²⁴等协议。
- WeDPR: WeDPR²⁵是微众银行推出的即时可用场景式隐私保护解决方案，依托区块链等分布式可信账本技术以及密码学技术，针对隐匿支付、匿名竞拍、匿名投票和选择性披露等应用场景落实用户数据和商业数据的隐私保护。

A.3 结合了 TEE 的区块链系统

尽管在上述部分方案中也采用了区块链或者分布式账本技术，但其主要目标仍然是提供通用的隐私保护方案。值得注意的是，还有一些项目定位于“具有隐私保护特性的区块链系统”，例如：

- TEEX: TEEX²⁶是一个基于公链的二层网络，通过将区块链链上执行过程下放到链下的可信执行环境中去，将执行过程和共识验算过程完全解耦开，从而减少公链的计算开销。
- Oasis: Oasis²⁷是由加州大学伯克利分校的 Dawn Song 教授领衔创办的区块链项目，旨在解决当下区块链在性能、安全、和隐私上的痛点。其特点在利用了 TEE 技术构建了智能合约执行平台 Ekiden，计算节点使用 TEE 来执行智能合约从而减少系统的共识开销。

可以发现，上述项目的核心思想是将原有链上执行逻辑（例如智能合约）放置于链下，从而减少共识的开销，同时基于 TEE 保证链下执行的可验证性。虽然涉及技术与其他隐私保护方案以及本项目存在相关性，但其应用场景有着较大区别，因此本文不做详细介绍。

²³https://en.wikipedia.org/wiki/Non-interactive_zero-knowledge_proof

²⁴<https://docs.ethhub.io/ethereum-roadmap/layer-2-scaling/zk-starks/>

²⁵<https://fintech.webank.com/wedpr/>

²⁶<https://teex.io/>

²⁷<https://www.oasislabs.com/>

附录 B 为什么隐私保护需要区块链

B.1 什么是区块链

区块链技术最初是被应用在比特币 [?] 系统中，随后在密码货币领域得到快速发展，并迅速扩展到物联网、知识产权保护等领域，被认为是继移动互联网之后的第五代互联网颠覆性技术。

区块链并不是一种单一的技术，而是多种技术整合的结果，包括但不限于分布式存储、共识机制、智能合约、对称/不对称加密等等。这些技术以新的结构组合在一起，形成了一种新的数据记录、存储和表达的方式。

从一般结构出发，区块链满足的基本特点有：

- 自治（透明）：系统节点对等，可以自由选择加入或离开；去中心化，无管理机构或第三方仲裁。
- 分布式（共享）：只需要连接到最近的节点就可以获取所有账本信息，同样交易发布也仅需提交给临近节点，依次转发下去直至遍及整个网络。
- 不可篡改（永久）：基于密码学技术保证上链交易合法，每个节点都本地同步一份附有时间戳的账本副本，保证数据不可篡改。
- 按合约执行（公平）：所有节点都按照一个规则或合约行事并达到共识（例如智能合约或账本同步）。

B.2 Fidelius 中的区块链做了什么

在区块链中，所有数据（包括交易类别、交易双方的地址、交易金额等）都是公开的，这在一定程度上提高了参与者对数据真实可靠的信心。大多数区块链系统会将所有交易数据记录在公共账本中，任何用户均可查询。一笔交易的有效性需要经过区块链中大部分节点的认可，而验证结果取决于交易数据（例如交易金额，携带数据及其签名）。

早期区块链系统（以公链为代表），依赖于公开透明的交易，所有节点同步账本，这种方式导致了严重的隐私泄露。此后产生了许多具有匿名区块链系统，例如 ZCash²⁸、Monero²⁹等等，其主要采用了混淆服务和环签名等技术。上述匿名技术主要是用于保护用户身份隐私，这一方面是由于公链属于非许可区块链，用户对于匿名

²⁸ZCash, <https://z.cash/zh/>

²⁹Monero, <https://www.getmonero.org/>

的需求更强烈；另一方面，大部分公链所承载的应用仍然为数字货币交易，即链上行为多数为普通转账交易，在此基础上，用户身份匿名化即可满足隐私保护需求。

然而联盟链的应用场景更为广泛，其链上主要数据不再是代币（Token）转账信息，更多的是业务相关数据，此时公开所有数据在某些关键应用场景下是不可行的，例如金融、医疗等数据敏感性高的场合。因此，区块链本身并不能直接满足交易数据的隐私保护需求。

近年来也涌现出一些隐私保护技术，可以对区块链的更多交易信息（例如智能合约）实现隐藏，例如在附录 A 中我们介绍的 zk-SNARK、WeDPR 等技术。然而在 Fidelius 的设计哲学中，我们并没有将 Fidelius 设计成与区块链紧耦合的关系，这是因为：

- Fidelius 的初衷是为了解决数据合作的隐私保护问题，而不仅仅局限于区块链上的数据隐匿；
- 我们发现当前大部分企业很难实现业务全部上链，这一方面是由于区块链作为底层基础设施仍然存在许多缺陷尚待解决³⁰，另一方面企业往往已经有成熟的数据存储平台（数据中心或者云存储），没有必要再基于区块链进行数据存储。

尽管如此，Fidelius 中也使用了区块链技术，这里区块链主要起到了如下作用：

- 可信第三方：如图 2 所示，数据提供方和数据使用方之间基于区块链实现了交互，由于共识机制本身的去中心化特性，实际上区块链扮演了可信第三方的作用。具体来说，区块链能够 1) 存储数据，且数据不会被篡改，2) 传输数据，且不依赖于中心化的证书机制（例如 HTTPS 所依赖的证书）³¹，3) 验证数据，即对数据使用情况及分析结果的验证。
- 抵抗作弊：由于区块链的特点，被区块链记录的数据合作行为是公开、且不可篡改的，这对于试图通过“刷单”等来影响数据合作的作弊行为有非常好的抵抗作用。虽然不同于在中心化的数据合作平台上处理作弊行为的“黑盒子”，公开的、不可篡改的数据合作记录可以被不断更新的自动化算法或人力检查，是一种有效的抵抗作弊的手段。

³⁰ 例如受共识机制、节点规模、存储等影响，区块链的吞吐量和扩展性仍然不能满足企业级业务的需求

³¹ 并且在参与主体规模增加的情况下区块链对传输开销的增加也非常低，若采用专线传输，开销则是指数级别的增长