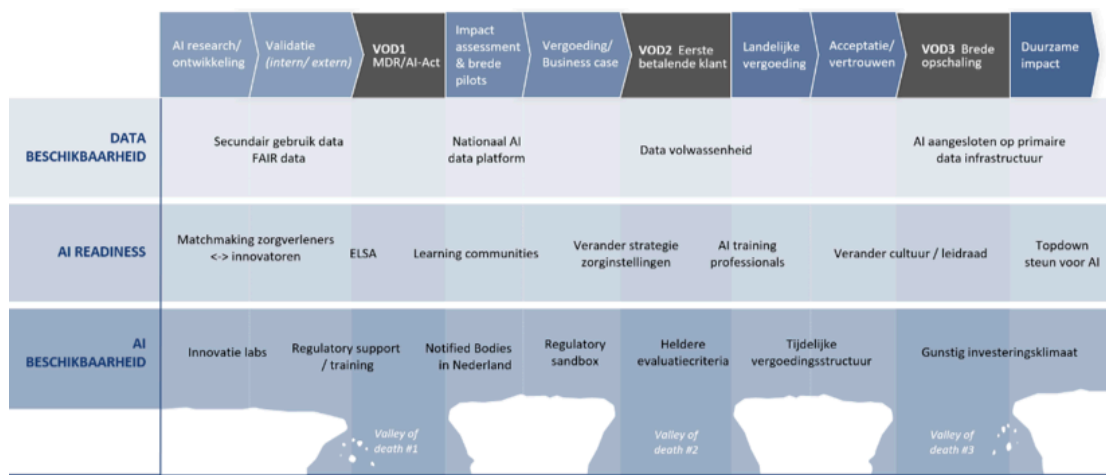


# Van terpen naar deltawerken voor AI in de zorg

## Startnotitie 'AI dataplaformen' (werktitel: AIDA)

### 1 Aanleiding: oproep voor een nationaal actieplan AI4health

Tijdens de laatste ICT&Health conferentie is een oproep gedaan om te komen tot een nationaal actieplan AI4Health. De kern van deze oproep is dat, ondanks alle lopende initiatieven, er nog steeds veel barrières zijn om data- en AI-gedreven innovaties in de zorg op grote schaal te realiseren. De praktijk wijst uit dat er drie *valleys of death* overkomen moeten worden (zie Gude W, Eekeren P van, Vasseur J [1] voor een recent overzicht):



Figuur 1: De drie *valleys of death* die grootschalig gebruik van AI in de zorg in de weg staan.

- Van concept tot toegang tot de praktijk.** Voordat grootschalig pilots en marktanalyse mogelijk is moet voldaan worden aan de voorwaarden van de Medical Device Regulation (MDR) en de AI act; 80-90% van de innovaties strandt hier. In de verschillende meetings en workshops die NLAIC georganiseerd heeft met zorginstellingen en innovatoren kwam de MDR/AI ACT bij herhaling als grote bottleneck naar boven (altijd top-3).
- Van pilot naar eerste betalende klant.** Dit vereist onder meer een solide business case, betrouwbare toegang tot data en professionele ondersteuning. Dit blijft een hardnekkig probleem, niet alleen voor toepassingen in de ziekenhuiszorg maar komt ook terug bij GGZ en VVT. Binnen de ziekenhuiszorg is het adresseren van deze VoD daarom bv een speerpunt voor de SAZ-ziekenhuizen (Expertisecentrum Zorgalgoritmen).
- Van eerste klant naar duurzame opschaling.** Implementatie, acceptatie en validatie ("calibratie") in andere omgevingen dan die van de pilot sites en eerste klant zijn verre van triviaal voor AI-innovaties. De meeste innovaties die de eerste klant weten te bereiken stranden alsnog in deze Valley of Death. Ter verdere illustratie van de

huidige stand van zaken: bij een AI-readiness traject van NLAIC kwam naar voren dat nog geen enkele zorginstelling daadwerkelijk AI-ready is.

## 2 Behoeftte aan een “AI dataplatform” (werktitel: AIDA)

Op dit moment wordt gewerkt aan een nationaal actieplan om deze barrières te adresseren langs vier actielijnen, te weten

1. Databeschikbaarheid
2. AI-readiness
3. AI-beschikbaarheid
4. Orkestratie.

Binnen de actielijn databeschikbaarheid is in de afgelopen jaren het nodige in gang gezet, waaronder het systematisch toepassen van FAIR principes, een helder en breed gedragen ethisch en juridisch kader en het realiseren van gedistribueerde toegang tot data. Tegelijkertijd kunnen we constateren dat het ontbreekt aan voldoende gedetailleerde afspraken om te komen tot een ecosysteem van ‘AI dataplatformen’ waarop effectief onderzoek kan worden gedaan, algoritmes kunnen worden ontwikkeld etc. Een dergelijk platform is een geïntegreerd systeem dat AI-ontwikkelaars ondersteunt met toegang tot data, modellen en andere hulpmiddelen om AI-projecten te ontwikkelen en te verbeteren. Dit platform biedt toegang tot essentiële bronnen, zoals datasets voor het trainen van AI-algoritmen, basis AI-modellen en diensten zoals een ELSA-desk voor ethische en juridische vraagstukken.

Dergelijke ‘AI dataplatformen’ zijn in feite specifieke vormen van beveiligde verwerkingsomgevingen zoals in de datagovernance verordening artikel 2 lid 20 is gedefinieerd:

“beveiligde verwerkingsomgeving”: de fysieke of virtuele omgeving en organisatorische middelen om te zorgen voor de naleving van het Unierecht, zoals Verordening (EU) 2016/679 (de Algemene verordening gegevensbescherming), met name wat betreft de rechten van datasubjecten, intellectuele-eigendomsrechten, en handels- en statistisch geheim, integriteit en toegankelijkheid, alsook van het toepasselijke nationale recht, en om de entiteit die de beveiligde verwerkingsomgeving biedt in staat te stellen alle gegevensverwerkingsactiviteiten te bepalen en er toezicht op te houden, met inbegrip van het tonen, opslaan, downloaden en exporteren van gegevens en het berekenen van afgeleide gegevens door middel van computeralgoritmen;

Onder de werktitel “AIDA” willen we in de komende periode met experts, belanghebbenden en veldpartijen te komen tot een gecoördineerde realisatie van een dergelijke nutsvoorzieningen. Deze startnotitie en website is bedoeld als interactief discussie document, ter ondersteuning van dit consultatieproces.

Op dit moment zijn er ontzettend veel ontwikkelingen gaande die relevant zijn voor AIDA. In het onderstaande geven we een samenvatting van relevante initiatieven, waarna we een eerste scoping presenteren en vragen formuleren als start voor de discussie.

## 3 Context en relevante initiatieven

### 3.1 Europese initiatieven

#### 3.1.a Simpl Open

Het Europese [Simpl](#) is een “... *is an open source, smart and secure middleware platform that supports data access and interoperability among European data spaces.*” In januari 2025 zijn gedetailleerde architecturen en functionele beschrijvingen van **Simpl-Open** opgeleverd om, zijnde een open-source software stack waarmee we deze generieke integratie laag op een gestandaardiseerde manier willen realiseren.

De Simpl-Open architectuur is een gedetailleerde uitwerking van bestaande referentie architecturen en is compatible met de Data Spaces Support Center (DSSC) Blueprint ([versie 1.5](#)) en de International Data Spaces Reference Architecture Model (IDS-RAM) ([huidige versie 4](#), [draft versie 5](#)).

#### 3.1.b AI Factories

Vanuit de EU wordt ingezet op de realisatie van [AI Factories](#), zijnde faciliteiten “... *that leverage the supercomputing capacity of the EuroHPC Joint Undertaking to develop trustworthy cutting-edge generative AI models.*” Dit initiatief zit meer in de hoek van High Performance Computing (HPC), en wordt ook getrokken vanuit de EuroHPC Joint Undertaking om betrouwbare, *state-of-the-art* generatieve AI modellen te ontwikkelen.

SURF is op dit moment penvoerder om namens Nederland een aanvraag in te dienen om een [grootschalige Nederlandse AI-faciliteit](#) te realiseren.

#### 3.1.c InvestAI

Europa heeft op 11 februari het [InvestAI-initiatief](#) aangekondigd om 200 miljard euro aan investeringen te mobiliseren. Dit initiatief is o.a. gevoed door CAIRNE, de *Confederation of Laboratories for Artificial Intelligence Research in Europe* dat al langer pleit voor een [CERN voor AI](#). Op dit moment is het nog onduidelijk wat deze initiatieven concreet zullen betekenen voor AIDA.

#### 3.1.d TEHDAS2

[TEHDAS2](#) is een “... *joint action [that] prepares the ground for the harmonised implementation of the secondary use of health data in the European Health Data Space – EHDS.*” Het is een Europees, zorg-specifiek programma, en veel van de werkpakketten zijn direct relevant voor AIDA. Een van de zaken die nader uitgezocht moeten worden is hoe de generieke architectuur van Simpl Open (sector onafhankelijk) zich verhouden tot de ontwerpprincipes en keuzes die binnen TEHDAS2 zijn gemaakt.

## 3.2 Nederlandse context

### 3.2.a Twiin

[Twiin](#) is een samenwerkingsverband waarin zorgaanbieders, leveranciers en partners werken aan het Twiin Afsprakenstelsel. Dit [Afsprakenstelsel](#) omvat gedetailleerde uitwerkingen over alle lagen van de architectuur voor het beschikbaar maken van gezondheidsgegevens. Zo zijn duidelijke keuzes gemaakt om bijvoorbeeld te werken met FHIR-gebaseerde *notified pull*, het gebruik van BSN voor identificatie, het gebruik van eIDAS voor authenticatie en OAuth2 voor autorisatie.

### 3.2.b NUTS

[NUTS](#) ontwikkelt en beheert een nutsvoorziening die het delen van zorg-gerelateerde informatie over het Internet mogelijk maakt op een vertrouwelijke, veilige en toegankelijke manier. Het Nuts-netwerk maakt gebruik van internationale standaarden om een vertrouwenslaag op het Internet te realiseren. Die standaarden zijn geïmplementeerd in de Nuts-node: Open Source software die zonder licentiekosten door elke IT leverancier gebruikt kan worden. Leveranciers mogen er ook voor kiezen om zelf de standaarden te implementeren.

### 3.2.c Health RI nodes

*Last but certainly not least* hebben de [Health RI nodes](#) in de afgelopen jaren het nodige ontwikkeld. Binnen de nodes is gewerkt aan verschillende oplossingen en aandachtsgebieden, waaronder het [myDRE Trusted Research Environment](#), het [molgenis](#) data platform gericht op wetenschappelijk onderzoek en bioinformatica, het [BBMRI-NL beeldanalyse platform](#) om er een paar te noemen.

## 4 Vraagstelling AIDA

Op dit moment worden AI-gedreven innovaties vaak ontwikkeld op niet gestandaardiseerde infrastructuur, wat leidt tot zeer hoge of zelfs onbetaalbare kosten voor een willekeurig onderzoek of innovatie. In analogie zou je kunnen zeggen dat elk onderzoek, elke innovatie zijn eigen terp moet bouwen, alvorens het project daadwerkelijk kan starten. Met AIDA willen we toe naar een Deltawerken voor AI4health. AIDA is hierbij het Deltaplan: het ontwerp van gestandaardiseerde en interoperabele “AI dataplatformen” waarmee komende jaren generieke, landelijk dekkende voorzieningen gerealiseerd kunnen worden. In deze analogie zijn de Deltawerken de realisatie hiervan, waarbij het perspectief is dat we niet naar één AI platform streven, maar een ecosysteem van platformen die interoperabel zijn. Net zoals dat de Deltawerken een ecosysteem van dijken, waterkeringen etc. zijn die elkaar ondersteunen, elk met een eigen functie.

In de komende maanden willen we dus een aanzet geven tot het opstellen van het Deltaplan. Belangrijke uitgangspunten en vragen hierbij zijn (indicatief, niet bedoeld als volledige opsomming):

- Hoe komen we tot harmonisatie, en waar nodig standaardisatie van verschillende oplossingsrichtingen op maximale interoperabiliteit te realiseren.

- Hoe kunnen we bestaande en in ontwikkeling zijnde architecturen combineren tot een consistent plan?
- Hoe kunnen we middels generieke functies het vertrouwensmodel goed te implementeren?
- Hoe kunnen we bestaande initiatieven maximale ruimte geven om zelf te blijven doorontwikkelen, met daarbij een minimale set van afspraken

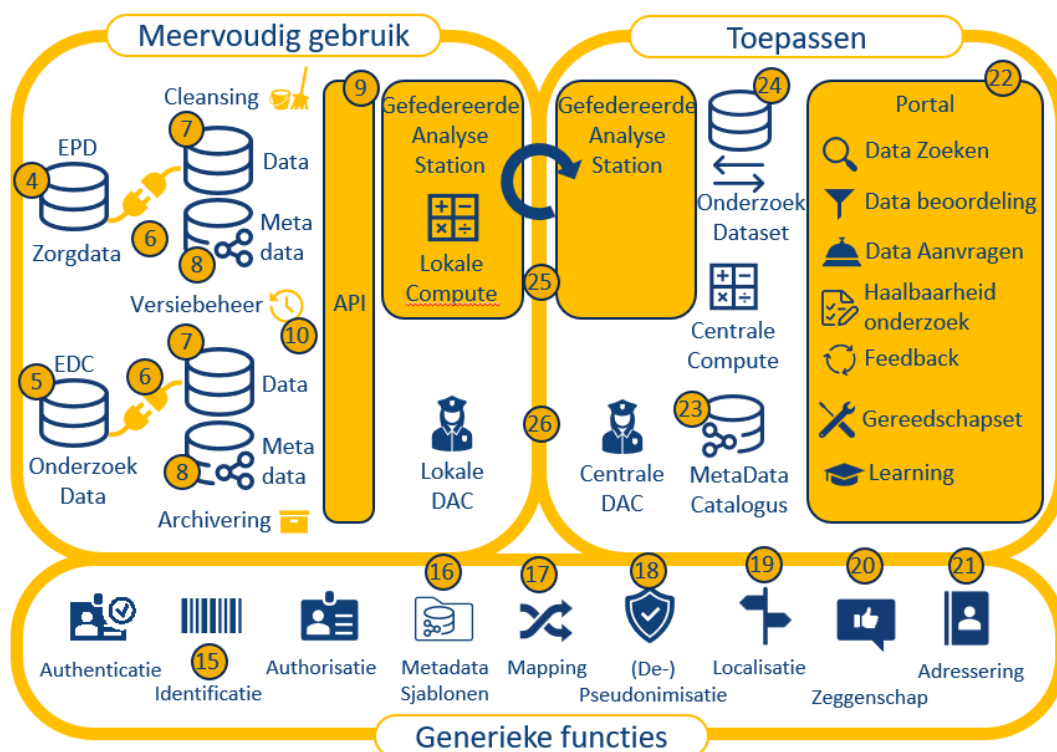
Bij het opstellen van deze startnotitie kwam vooral ook naar voren dat er behoefte is aan een eenduidige terminologie en beschrijvingen van componenten. Dat is een van de concrete *deliverables* van AIDA voor dit jaar.

Als start voor de discussie, benoemen we de volgende drie thema's die we met het veld verder willen verkennen en verdiepen:

1. Onderscheid verschillende vormen van *Secure Processing Environments* (SPEs)
2. Koppelvlak tussen data en SPEs
3. Orchestratie van infrastructuur

## 5 Thema 1: soorten SPEs

Versie 4 van de Health-RI wiki beschrijft de [gezondheidsdata-infrastructuur voor onderzoek, beleid en innovatie](#). Deze infrastructuur is specifiek gericht op secundair gebruik, en is een verbijzondering van de algemene [gezamenlijk gezondheidsdata architectuurmodel](#). Binnen deze architectuur zijn reeds twee soorten van Secure Processing Environments benoemd, namelijk [veilige verwerkingsomgevingen](#)) en [gefedereerde verwerkingsomgevingen](#).



Figuur 2: Conceptuele architectuur voor secundair gebruik

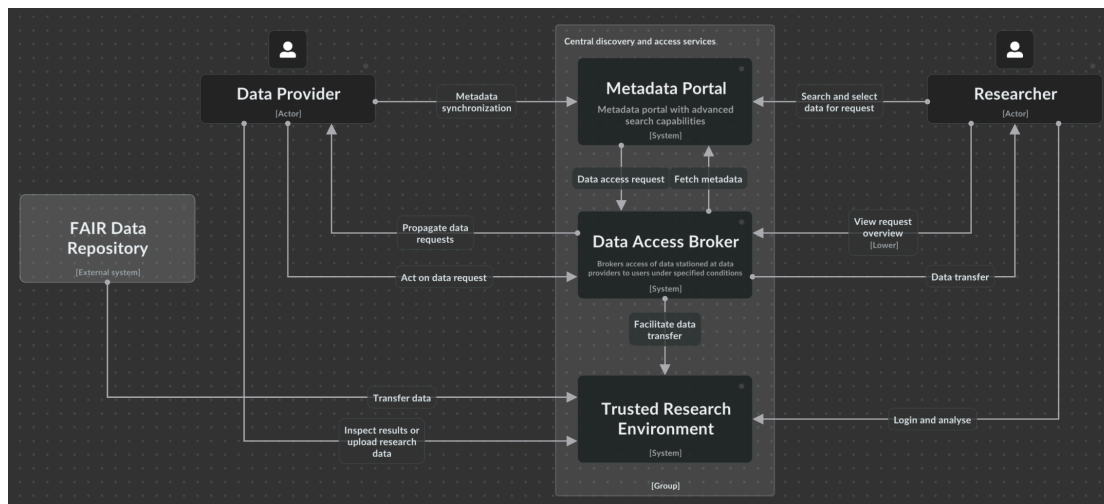
De basis gedachte achter AIDA is dat er *verschillende soorten van secure processing environments (SPEs)* zullen zijn. Daarbij introduceren we een derde type, de hybride SPE. In dit document hanteren we de namen van deze drie typen om expliciet onderscheid te maken; we zullen spreken van SPEs in het algemeen als we alle drie de soorten bedoelen. We geven een korte schets en voorbeelden van elk type.

Tabel 1: Drie soorten van Secure Processing Environments die binnen de scope van AIDA vallen.

Centrale SPE	<ul style="list-style-type: none"> <li>• vaak benoemd als Trusted Research Environment</li> <li>• veel bestaande voorbeelden, zie o.a. EOSC-ENTRUST</li> <li>• machine learning op tabulaire data mogelijk</li> <li>• <i>statistical disclosure control</i> op output</li> </ul>
Decentrale SPEs	<ul style="list-style-type: none"> <li>• decentrale c.q. federated benadering</li> <li>• oorspronkelijk bedoelt voor machine learning</li> <li>• kan ook gebruikt worden voor statistische analyse</li> <li>• moeilijker om mee te werken</li> </ul>
Hybride SPE (H-SPE)	<ul style="list-style-type: none"> <li>• Combinatie van bovenstaande technieken</li> <li>• Nodig om gebruik te maken centrale rekencapaciteit</li> <li>• Gedachte om gebruiksgemak te verbeteren</li> </ul>

## 5.1 Centrale: SURF Secure Analysis Environment (SANE)

SURF Secure ANALysis Environment (SANE) is een virtuele, volledig afgeschermdde omgeving waarop met vooraf goedgekeurde analyse software draait en toegang tot sensitive data wordt gegeven (Figuur 3). In onderstaand overzicht is SANE gepositioneerd als TRE, waarmee de data aanbieder controle houdt over de data die ter beschikking wordt gesteld en waarmee de data consumer op een makkelijke manier toegang krijgt. SANE biedt functionaliteiten op het gebied van *Research Analytics*, *Secure Data Zone* en *Data Discovery*. Meer details staan in de [blauwdruk van EOSC-ENTRUST](#).



Figuur 3: Positionering van SANE binnen een generieke data space architectuur.

Belangrijk kenmerk van SANE en andere TREs is dat de data fysiek naar de *Secure Data Zone* wordt gekopieerd. Naast het veilig aanbieden van data (als data houder), is dit ook het mechanisme waarmee data gebruikers hun eigen data mee kunnen nemen naar de TRE, om daarbinnen te koppelen aan andere data. Dit gebeurt vaak met gebruik van pseudonimisering. De CBS microdata omgeving werkt op een vergelijkbare manier.

Binnen de blauwdruk van EOSC-ENTRUST wordt gesproken over *Federation Services* tussen verschillende TREs. Daarbij gaat het om data federation: data wordt (tijdelijk) van de ene naar de andere TRE gekopieerd zodat het daar in combinatie verwerkt kan worden. Data federation als mechanisme is anders dan federated learning: daarbij worden de berekeningen decentraal uitgevoerd en alleen de resultaten centraal gedeeld (zie hieronder). Federated learning is met name nuttig voor horizontaal gepartitioneerde data. Voor verticaal gepartitioneerde data, is data federation zoals beschreven door EOSC-ENTRUST meer geschikt.

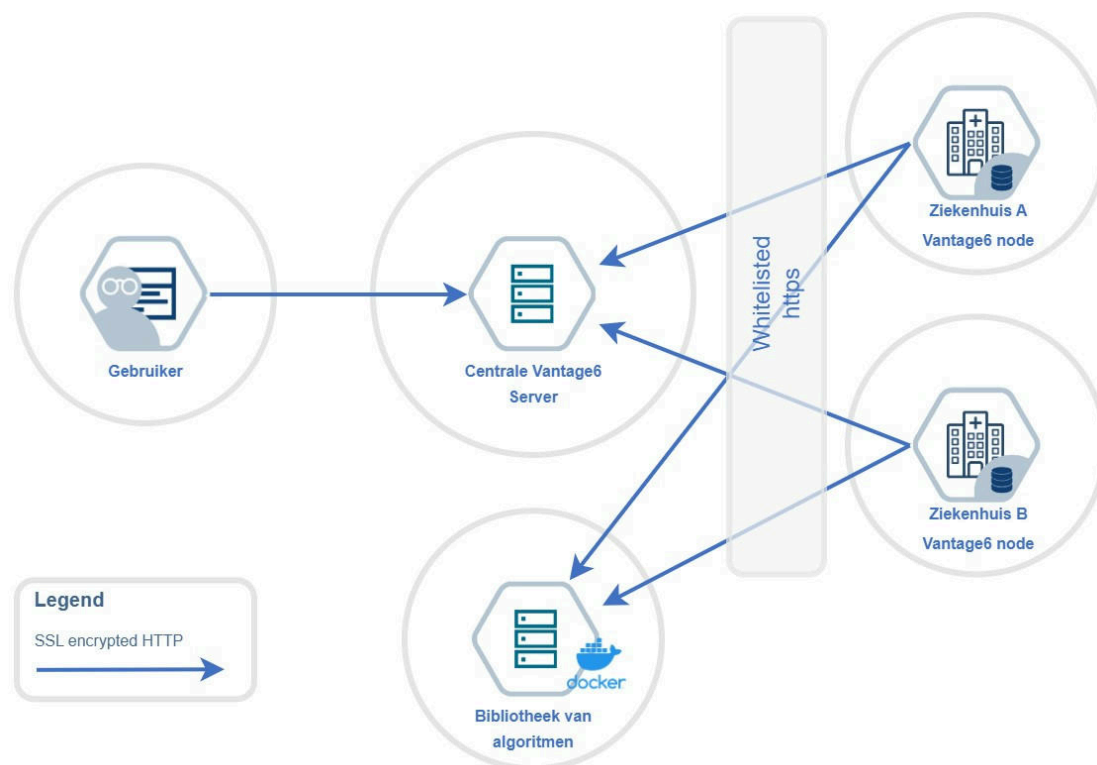
Er zijn meer voorbeelden van centrale SPEs. Zo hebben de meeste *National Statistics Offices* (NSOs) zoals het CBS een [microdata omgeving](#). Alhoewel deze omgevingen zijn opzet voordat machine learning zijn intrede deed, bieden de meeste microdata omgevingen nu ook al de mogelijkheid om 'lichte' algoritmes te trainen op tabulaire data. Een rapport van de Verenigde Naties beschrijft dat deze omgevingen in toenemende mate ook worden uitgebreid met nieuwe AI-technieken, zoals privacy-enhancing technologieën (PETs, zie [2]).

Er zijn ook voorbeelden van centrale SPEs specifiek voor de zorg:

- Het Finse Social and Health Data Permit Authority (Findata) biedt met [Kapseli](#) een landelijke voorziening aan dat aanvullend is op het Finse NSO.
- Het [Mayo Clinic Platform\\_Discover](#) is een voorbeeld van een platform binnen een netwerk van zorg leveranciers.

## 5.2 Decentrale SPEs: PLUGIN/vantage6

Decentrale SPEs maken gebruik van federated learning (FL), wat als concept ook wel bekend staat als de Personal Health Train (PHT). FL wordt in toenemende mate gebruikt in de zorg [3]; de term FL wordt vooral gebruikt om naar het technische concept te verwijzen, terwijl PHT verder gaat in het definiëren van afspraken rondom het gebruik van FL. In Nederland is een actieve community rondom het [vantage6 platform](#) dat wordt gebruikt in het [PLUGIN project](#), en internationaal in [50 andere netwerken](#). Het basis principe is dat bij FL de gegevens op afzonderlijke 'data stations' verschillende apparaten blijven die participeren in het federatieve netwerk. Om deze data te gebruiken voor machine learning, wordt bij elk data station het algoritme lokaal c.q. afzonderlijk getraind. Vervolgens worden alleen de resultaten van het algoritme - bijvoorbeeld geaggregeerde statistieken of de modelparameters van het neurale netwerk - gedeeld met een centrale server. Deze server combineert de resultaten van afzonderlijke modellen tot één model, welke vervolgens met alle deelnemers van het federated SPE gedeeld wordt.



Figuur 4: Overzicht van vantage6 infrastructuur zoals in PLUGIN is gerealiseerd.

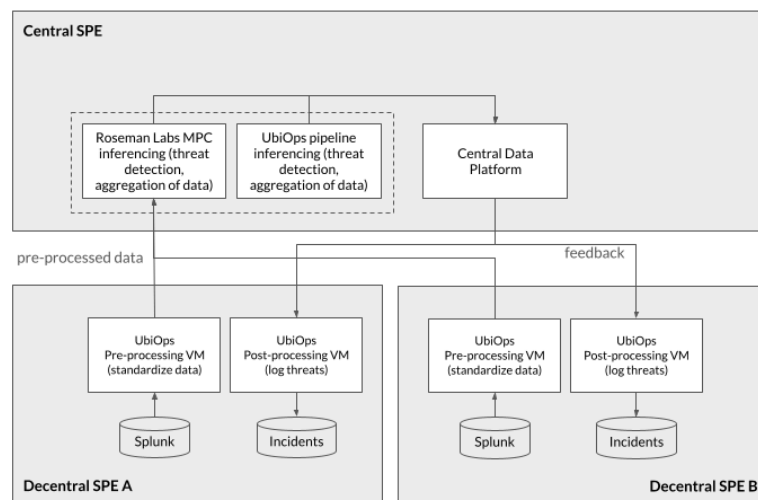


Het [PLUGIN project](#) heeft een decentrale SPE van tientallen ziekenhuizen gerealiseerd, waarbij gebruik wordt gemaakt [vantage6](#) als platform. Belangrijkste kenmerken van deze opzet zijn:

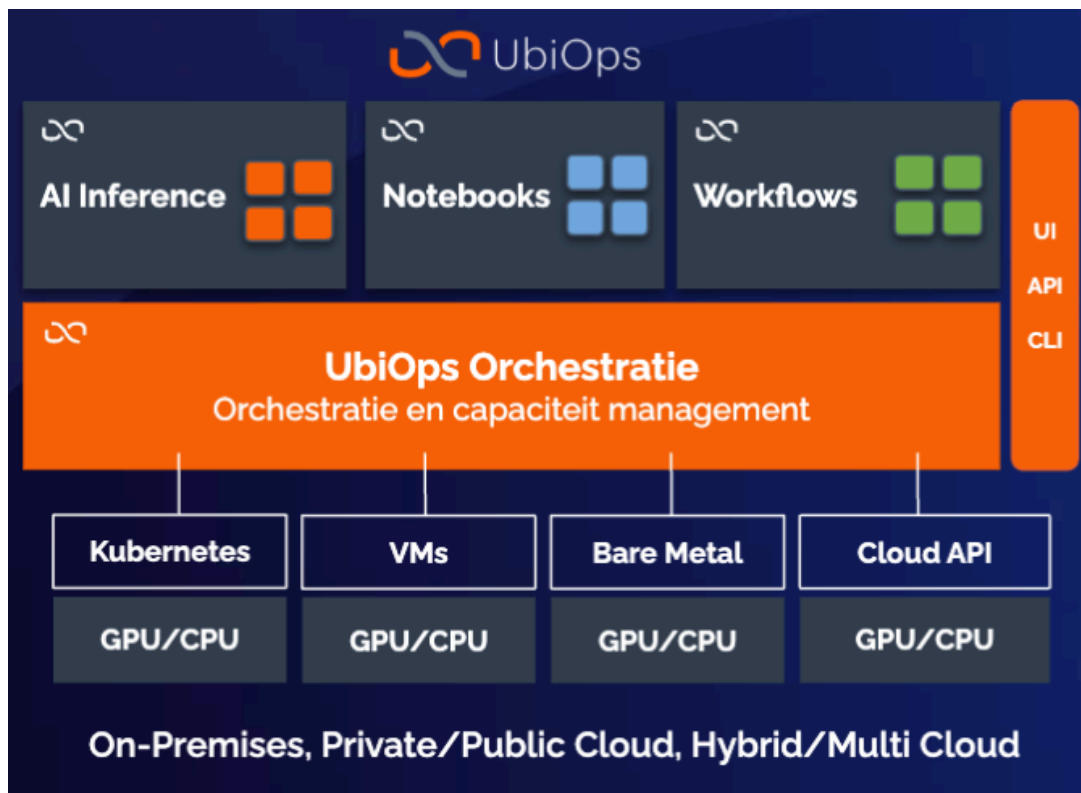
- Gebruik van beveiligde containers en virtual private networks voor de infrastructuur laag;
- Ontzorgen van deelnemende ziekenhuizen, waarbij gebruik wordt gemaakt van een generieke Linux server in het IT domein van het ziekenhuis dat als basis dient om de opslag en rekenkracht om het vantage6 netwerk te realiseren. Afhankelijk of een ziekenhuis mee doet als trainingsziekenhuis of alleen als gebruiker dient een zwaardere resp. lichtere Linux server te worden geconfigureerd;
- Voor elk project wordt de berekening c.q. machine learning expliciet ‘verpakt’ in een Docker container, zijnde de berekening die daadwerkelijk wordt uitgevoerd;
- De generieke Linux server wordt ook gebruikt om dashboard, informatieproducten etc. te hosten binnen het IT domein van het ziekenhuis.

Het gebruik van een standaard data model (op de data stations) is een belangrijke randvoorwaarde om FL te kunnen doen. Naast het gebruik van vantage6 als kerntechnologie, heeft PLUGIN ervoor gekozen om FHIR als data standaard te gebruiken. Hiertoe is een [FHIR profiel in ontwikkeling](#) die aansluit op de bestaande ZIBS2020 bouwstenen. Meer achtergrond over de keuze voor FHIR is te lezen in [dit artikel](#). Andere voorbeelden van decentrale SPEs gebaseerd op FL zijn [hier](#) te vinden.

### 5.3 Hybride SPE: UbiOps en Roseman Labs



Figuur 5: Hybride SPE in de veiligheidsketen tussen verschillende Security Operating Centra (SOCs).



Figuur 6: UbiOps orkestratielaag

De hybride SPE is een nieuwe oplossingsrichting die we willen verkennen en realiseren in AIDA. Er zijn minder concrete voorbeelden van een dergelijke opzet. UbiOps en Roseman Labs hebben een oplossing die er het dichtst bij in de buurt komt (Figuur 5). In analogie met data spaces, gaat het hier om het verbinden van verschillende Security Operating Centra (SOC) in de beveiligingsketen. In een hybride SPE kunnen *compute* (rekenkracht) en *storage* (opslag) zowel lokaal als centraal worden uitgevoerd. In deze architectuur worden bijvoorbeeld pre-processing van data decentraal uitgevoerd in de SOC's in de onderste laag van Figuur 5. De resultaten van deze pre-processing gaan naar het centrale platform. Daar kunnen vervolgens ook weer (vervolg-)berekeningen worden uitgevoerd, op de *storage* en *compute* die beschikbaar zijn in de centrale SPE.

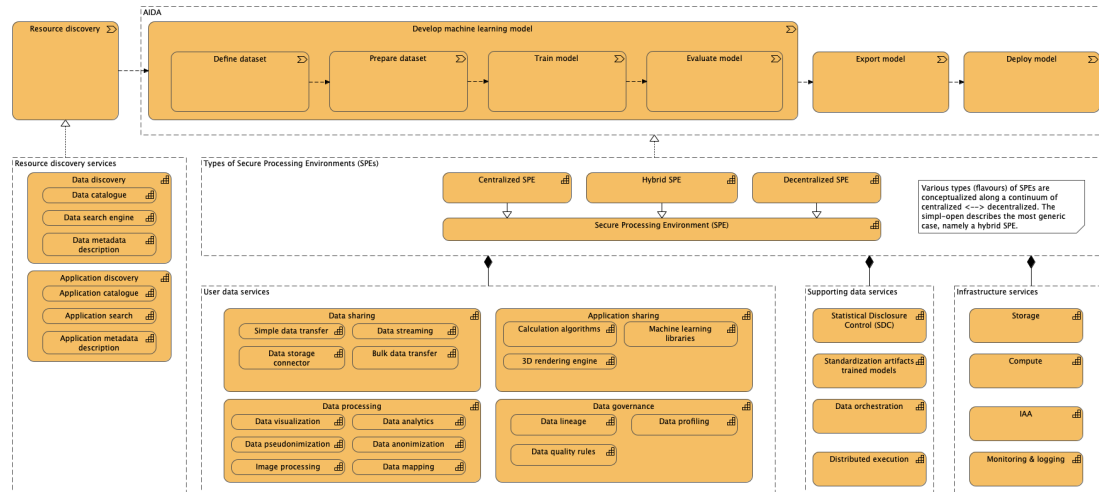
De hybride SPE is mogelijk gemaakt door [UbiOps](#), een platform leverancier die de orkestratielaag biedt waarmee alle *storage* en *compute* centraal wordt beheerd (Figuur 6). Een belangrijk ontwerpprincipe van deze orkestratielaag is dat het verschillende soort fysieke infrastructuur kan managen, variërend van *bare metal* servers, Kubernetes cluster, virtuele machines, public cloud infrastructuur etc.

Een ander onderscheidende kenmerk van deze opzet is dat de centrale dataverwerking ook onder encryptie uitgevoerd kan worden via het Roseman Labs MPC (Multiparty Computation) platform. Door berekeningen *in-the-blind* uit te voeren, zijn de data extra beschermd.

Deze opzet van hybride SPE is in lijn met de recent gepubliceerde [Simpl-Open architectuur](#). Deze aanpak biedt de mogelijkheid om over verschillende SPEs tot harmonisatie en

interoperabiliteit te komen. Denk bijvoorbeeld aan een situatie waarbij een analyse kan worden uitgevoerd over verschillende Health-RI nodes heen. In Paragraaf 7 gaan we hier dieper op in.

## 5.4 Strategy view op AIDA

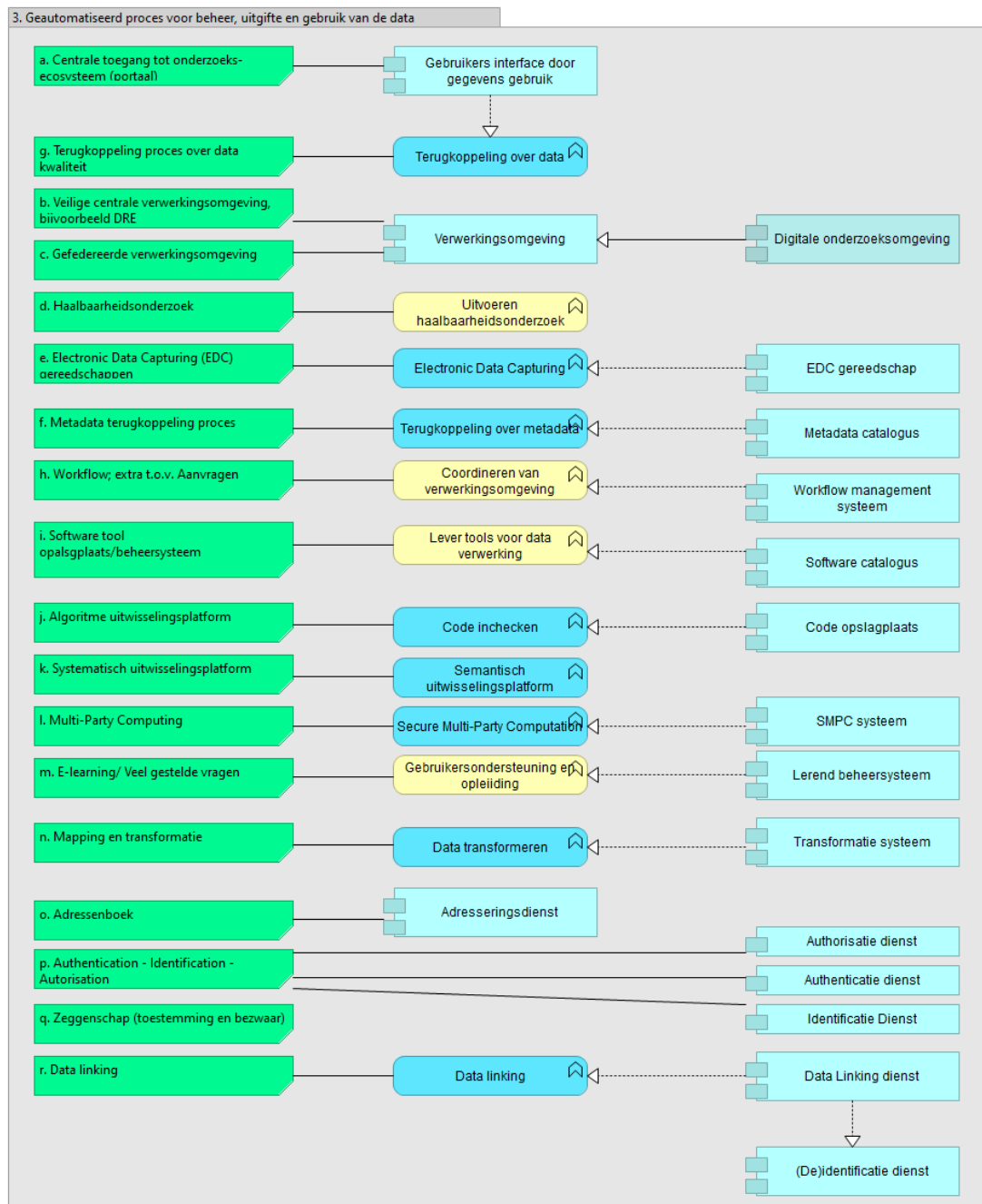


Figuur 7: De Strategy view als startpunt voor de discussie.)

Gegeven deze verschillende soorten SPEs is een eerste *strategy view* van AIDA geschetst in Figuur 7. De *value stream* elementen zijn beschreven in termen van het ontwikkelproces van **CRISP-DM**. Deze *value stream* kan worden gerealiseerd met behulp van verschillende soorten SPEs. De modulaire *capabilities* zijn de verschillende functionele bouwblokken die in een SPE gebundeld/aangeboden kunnen worden. De gedachte is dat elke SPE, afhankelijk van de context, doelgroep etc. een eigen configuratie van *capabilities* heeft.

## 6 Thema 2: koppelvlak datastores en SPE

Het ontwerp van AIDA vereist harmonisatie van de koppelvlakken (component 9 in Figuur 2) tussen de datastores (component 7) en de SPE. Binnen de Simpl Open architectuur wordt gesproken van Resource discovery (zie Figuur 7, linker cluster van *capabilities*). Versie 4 van de Health-RI wiki geeft een uitwerking van deze *discovery services* in een aantal bouwblokken (Figuur 8). In het onderstaande gaan we in op een aantal van deze bouwblokken die relevant zijn voor de uitwerking van AIDA



Figuur 8: Bouwblokken voor geautomatiseerd proces voor beheer, uitgifte en gebruik van de data (bron).

## 6.1 Onderscheid datastores en datasets

Een van de uitgangspunten van de Health-RI architectuur is dat data zo dicht mogelijk bij de bron FAIR wordt gemaakt. Het concept van een FAIR Data Point is hierin leidend, waarvan inmiddels verschillende implementaties zijn en complementaire software componenten voor b.v. het *harvesten* van FAIR metadata catalogi.

Een van de zaken die nader gespecificeerd moet worden is het onderscheid tussen een FAIR data set en een datastore. Een datastore in deze is een component dat als *system*

of record data beschikbaar stelt voor hergebruik en/of een longitudinale datastore waarbij data vanuit verschillende bronsystemen worden gecombineerd en gepersisteerd. openEHR en OMOP zijn veel gebruikte standaarden voor de implementatie van een dergelijk datastore. Voor meer details en achtergrond verwijzen we naar Tsafnat G, Dunscombe R, Gabriel D, Grieve G, Reich C [4].

Met name voor het beschikbaar stellen en hergebruiken van tabulaire data uit EPDs, HISen etc. is het wenselijk om snel te verkennen of een bepaalde dataset relevant is voor een gebruiker. We zien ‘verkenner’ voor ons, waarbij potentiële data gebruikers in staat worden gesteld om de datastore te queryen. Typisch betreft dit het zoeken naar relevante populaties in de datastores, bijvoorbeeld “hoeveel patiënten ouder dan 65 met aandoening XXX” zitten in de datastore. De resultaat van een dergelijk query is overigens altijd een anonieme dataset met enkel geaggregeerde, beschrijvende statistieken.

De meeste - en zo niet alle - open standaarden voor datastores hebben een mechanisme om dit te ondersteunen:

- openEHR gebruikt hiervoor het [Archetype Query Language \(AQL\)](#)
- OMOP gebruikt SQL (geen OMOP-specifieke query language) direct op de relationale database
- binnen het FHIR ecosysteem zijn nieuwe standaarden zoals [Bulk FHIR](#) en [SQL-on-FHIR](#) bedoelt om subsets te definiëren en op een makkelijke manier data te aggregeren

Deze benadering zorgt ervoor dat (potentiële) data gebruikers op een eenduidige manier ‘inclusiecriteria’ kunnen bepalen op de datastore. Het gaat hierbij om inclusie van rijen (subjecten, patient) en kolommen (welke attributen, waarden). Alhoewel het runnen van de query elke keer een iets ander resultaat geeft (je kunt meer rijen krijgen omdat nieuwe patiënten zijn toegevoegd), is het zodanig gestandaardiseerd dat het een vergelijkbare functie vervult als een statische, gepersisteerde dataset zoals een FDP dat beoogt. De metadatering van de query op de datastore zal waarschijnlijk met nagenoeg dezelfde metadata standaarden (DCAT-AP) kunnen worden gedaan.

Alhoewel de inrichting van een ‘verkenner’ functie technisch mogelijk is, blijkt uit de praktijk dat er nog het nodige aan data verificatie, mapping etc. gedaan moet worden om te komen tot een ‘bevroren’ dataset die geschikt is voor secundair gebruik. Tijdens de hackathon in december 2024 zijn een aantal knelpunten geïdentificeerd. In de uitwerkingen van AIDA willen we expliciet ingaan hoe we dit kunnen oplossen.

## **6.2 Standaard voor koppelen van verticaal gepartitioneerde data middels pseudonimisering**

In veel use-cases dient data uit verschillende bronnen te worden gekoppeld op subject/patient niveau. Op dit moment is ZorgTTP een veelgebruikte pseudonimiseringsmethode dat wordt gebruikt, dat sinds 2015 tevens [openbaar is](#). Het gebruik van pseudoniemen om te koppelen is met name relevant voor centrale en hybride SPEs: op deze manier kan een integrale dataset gemaakt worden voor secundair gebruik. Een decentrale SPE is

minder geschikt voor het koppelen van zogenaamde verticaal gepartitioneerde data; het is niet onmogelijk maar zeker bewerkelijker.

## 7 Thema 3: Orchestratie van infrastructuur

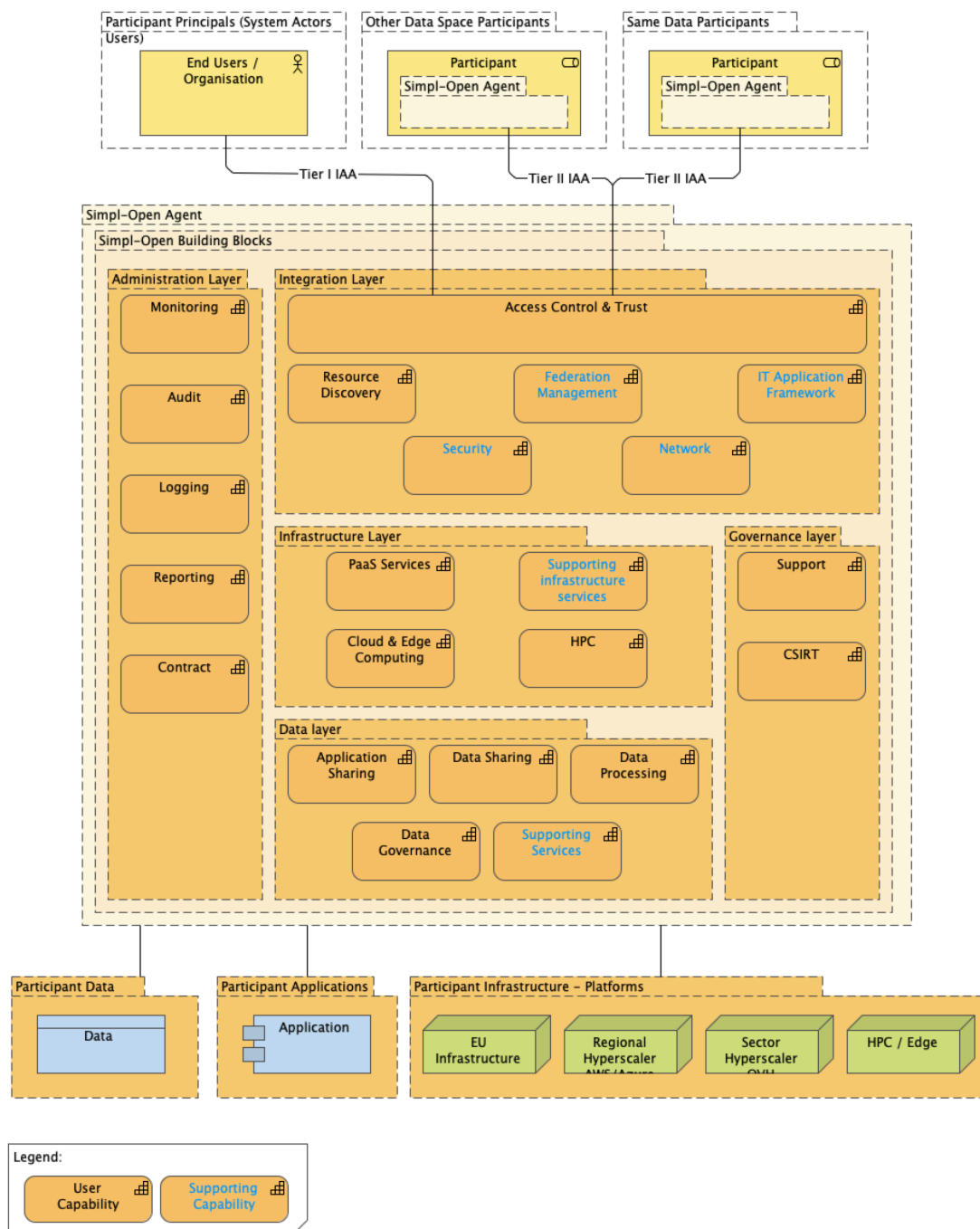
Orchestratie van infrastructuur binnen en tussen verschillende AIDA platformen is essentieel om tot een samenhangend Deltawerken te komen. Binnen de Simpl Open architectuur is de *Infrastructure Layer* ook in detail uitwerkt.

### **i** Infrastructure Layer in Simpl Open

The capabilities provided by this layer enable the consumers to easily provision the necessary computing and storage resources to execute their workloads in a secure and energy-efficient way. The **infrastructure orchestration**, automates the provisioning of the infrastructure resources to enable the various infrastructure providers to interconnect infrastructure orchestration and get exposed via a standard interface. The **distributed execution**, allows the consumer to deploy applications and execute computations close to the distributed execution data.

The **cloud & edge computing** capability provides the opportunity to provision various resources to execute computations or store data in the environment of their choice. The **platform-as-a-service** building blocks provide several database engines and other platform-level resources. Finally, the **HPC** capability permits the consumer to perform complex calculations at high speed by providing a cluster of high-performance computers.

The infrastructure building blocks can be easily combined with each other to create even more value for the consumer. For instance, after successfully analysing certain data with the help of the provisioned *platform-as-a-service* analytical resources or the *HPC* capability, the consumer may want to store the used datasets and/or the results of their calculations. In this case, the *PaaS* storage building block can easily fulfil the storage needs of the end user. In case a consumer would like to develop a stand-alone application, they may also use various *PaaS* resources at the same time. They can leverage the different storage options to store each sort of data in the most efficient manner (e.g. the transactional data in a transactional database, while the sensor data in a NoSQL database). Besides, they can use the *cloud & edge computing* capabilities to deploy and run their applications, and the *distributed execution* capability even enables them to run the code close to the edge.



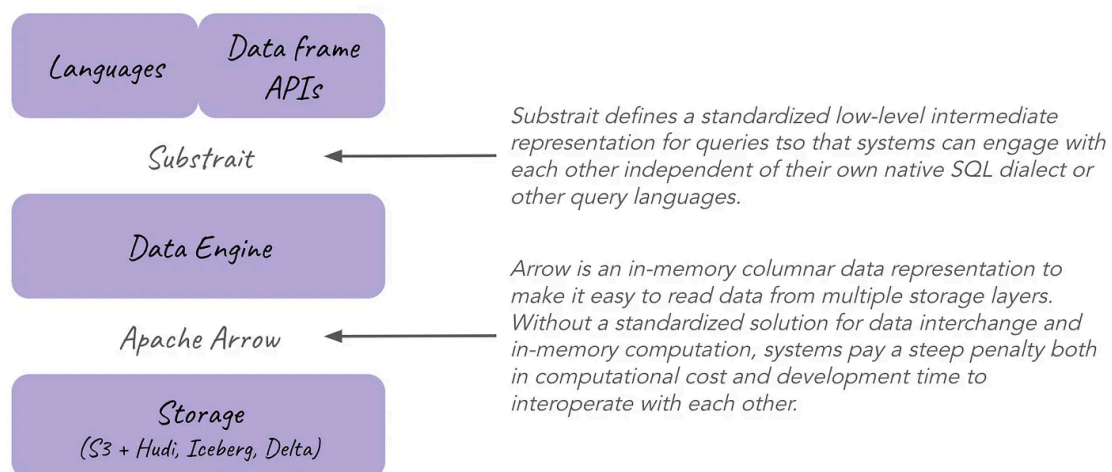
Figuur 9: Conceptueel overzicht Simpl Open, met een centrale plek van de *Infrastructure Layer*

Simpl Open hanteert de volgende ontwerpprincipes voor de infrastructuur laag:

- Het moet mogelijk zijn om over verschillende fysieke locaties een data space op te zetten
- Het gebruik van *agents* als mechanisme voor het orkestreren van allerlei *compute* en *storage* binnen een data space en tussen een data space

- Indeling van twee Tiers voor Identificatie, Authenticatie en Autorisatie (IAA) functies, waarvoor standaarden gebruikt moeten worden
  - Tier 1: IAA van gebruikers
  - Tier 2: IAA voor machine-to-machine orkestratie
- Naast gebruik van catalogi voor data en applicaties wordt ook het gebruik van een infrastructuur voorgeschreven, zodat daarmee inzichtelijk is welke *compute* en *storage* beschikbaar is binnen het netwerk
- Aansluiten bij bestaande standaarden en open source componenten zoals Keycloak als IAA applicatie component, OAuth2 voor autorisatie, S3-compatible blob storage als default opslag dienst en Kubernetes voor de containerinfrastructuur

De modulaire architectuur van Simpl Open heeft veel raakvlakken met de huidige trend binnen de data engineering community om toe te werken naar een zogenaamde *composable data stack*.<sup>1</sup> In deze *composable data stack* worden *storage*, *compute* en *query languages* volledige 'ontvlochten' in afzonderlijke componenten. Interoperabiliteit tussen de componenten is gebaseerd op moderne standaarden zoals Apache Arrow, Apache Parquet, Apache Iceberg en Apache Substrait (zie Figuur 10).



Figuur 10: Vereenvoudigd overzicht van de *composable data stack* (bron).

## Bibliografie

1. Gude W, Eekeren P van, Vasseur J (2024) AI Monitor Ziekenhuizen 2024
2. (2023) The PET Guide
3. Teo ZL, Jin L, Liu N, et al (2024) Federated Machine Learning in Healthcare: A Systematic Review on Clinical Applications and Technical Architecture. Cell Reports Medicine 5(2):101419. <https://doi.org/10.1016/j.xcrm.2024.101419>
4. Tsafnat G, Dunscombe R, Gabriel D, Grieve G, Reich C (2024) Converge or Collide? Making Sense of a Plethora of Open Data Standards in Health Care. Journal of Medical Internet Research 26(1):e55779. <https://doi.org/10.2196/55779>

<sup>1</sup>Zie de [Composable Codex](#) voor een uitgebreide toelichting.