

Health
Campus

Den
Haag

VIPP: Het Virtuele Patiënten en Populatie project

“Een digitale tweeling mét ELAN”

Symposium CBS, 1 juni 2023, Prof. Dr. Marco Spruit, LUMC/LIACS

“Synthetische data buiten het CBS en in het LUMC”



Agenda

1. Synthetische data in LUMC context? [Over mij, HCDH en ELAN](#)
2. Onderwijs met synthetische data? [PRIMA 2020](#)
3. Onderzoek met synthetische data? [Drie benaderingen](#)
4. Et cetera? [Effecten en de nabije toekomst](#)

**Health
Campus**

**Den
Haag**

Synthetische data in LUMC context

Ikzelf, de Health Campus Den Haag, ELAN datawarehouse



Mattijs Numans Jeroen Struijs Dennis Mook

Over... Marco Spruit

Pretpakket VWO
Propedeuse Nederlands
Propedeuse Muziekwetenschap
Bovenbouwstudie α -Informatica



Als Ontwikkelaar



- 1993
 - *Information Retrieval* programmeur
 - ZyLAB Europe BV



- 1995
 - *Big Data* systeemontwikkelaar
 - Koninklijke Marine



- 1997
 - Productsoftware ontwikkelaar/ondernemer
 - Insertable Objects, Wizzer BV

Als Wetenschapper



- 2003
 - OiO in Computacionele Linguïstiek
 - Universiteit van Amsterdam (UvA)



- 2007
 - Universitair (hoofd)docent Informatiekunde
 - Universiteit Utrecht >> ADS Lab (2015)

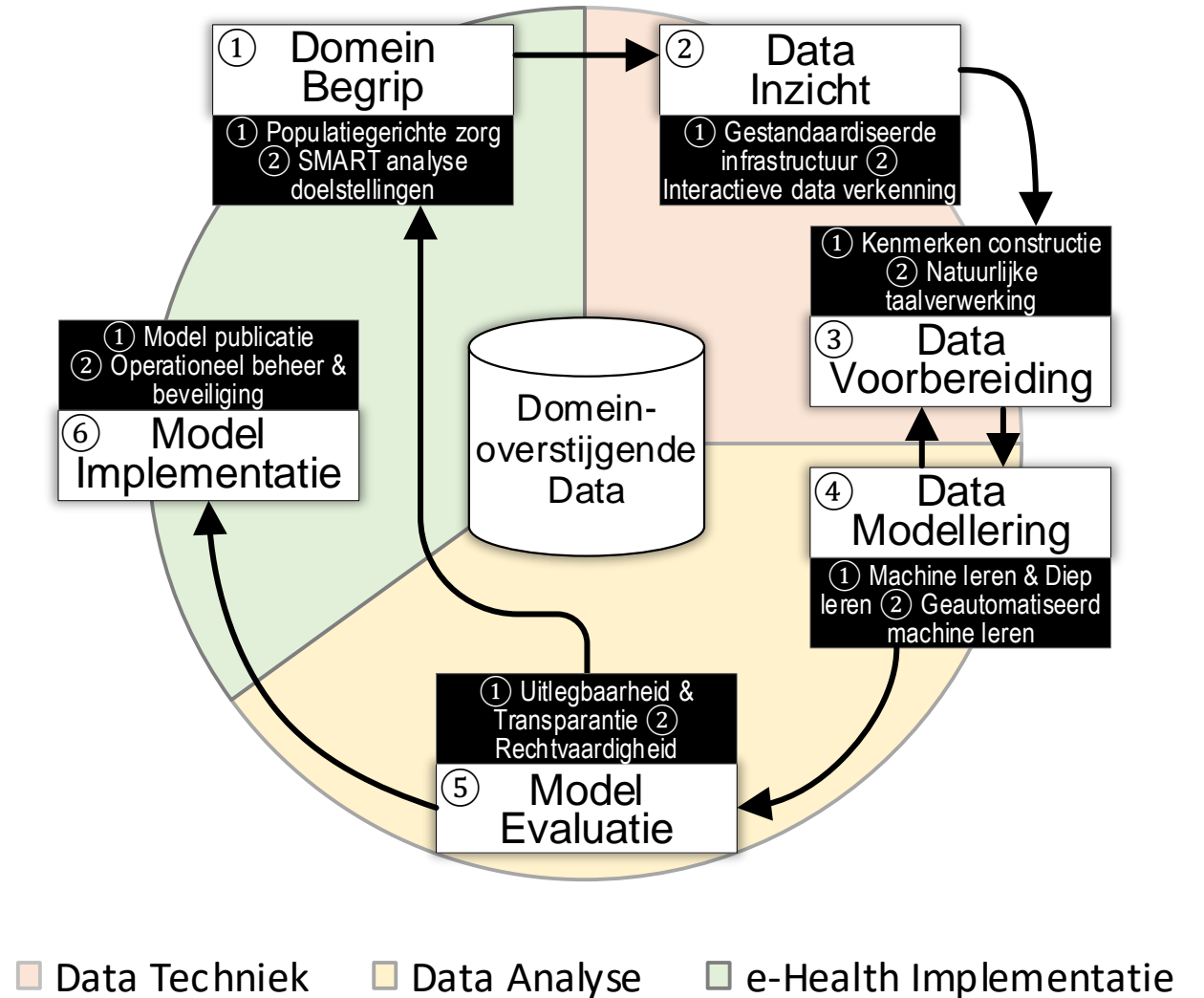
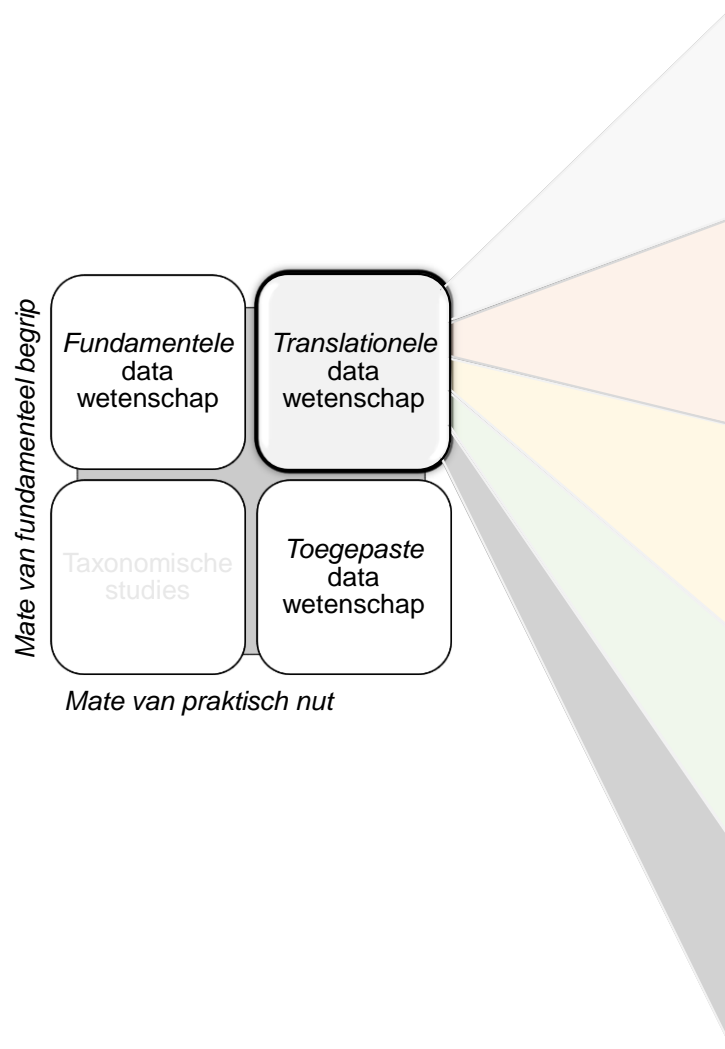


- 2020
 - **Hoogleraar *Advanced Data Science in Population Health***
 - **LUMC/LIACS @ Universiteit Leiden (HCDH)**
 - PH Living Lab, CAIRELab, TDS Lab, SIG HDS

Health
Campus

Den
Haag

Over... *Translationele* Datawetenschap



Over... Health Campus Den Haag

1. Verkleinen gezondheidsverschillen
2. Duurzame benadering
3. Breed gezondheidsperspectief



Populatiegericht
zorgonderzoek



Mentale gezondheid



Universiteit
Leiden

Data science &
AI



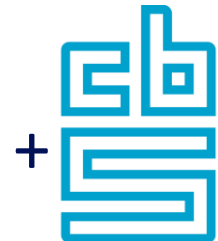
Zorg/Data-
onderzoek



Ziekenhuis data



Den Haag
oa WMO data



Ziekenhuis data



hadoks

Huisarts data



oa Vaccinatie data



Ziekenhuis data



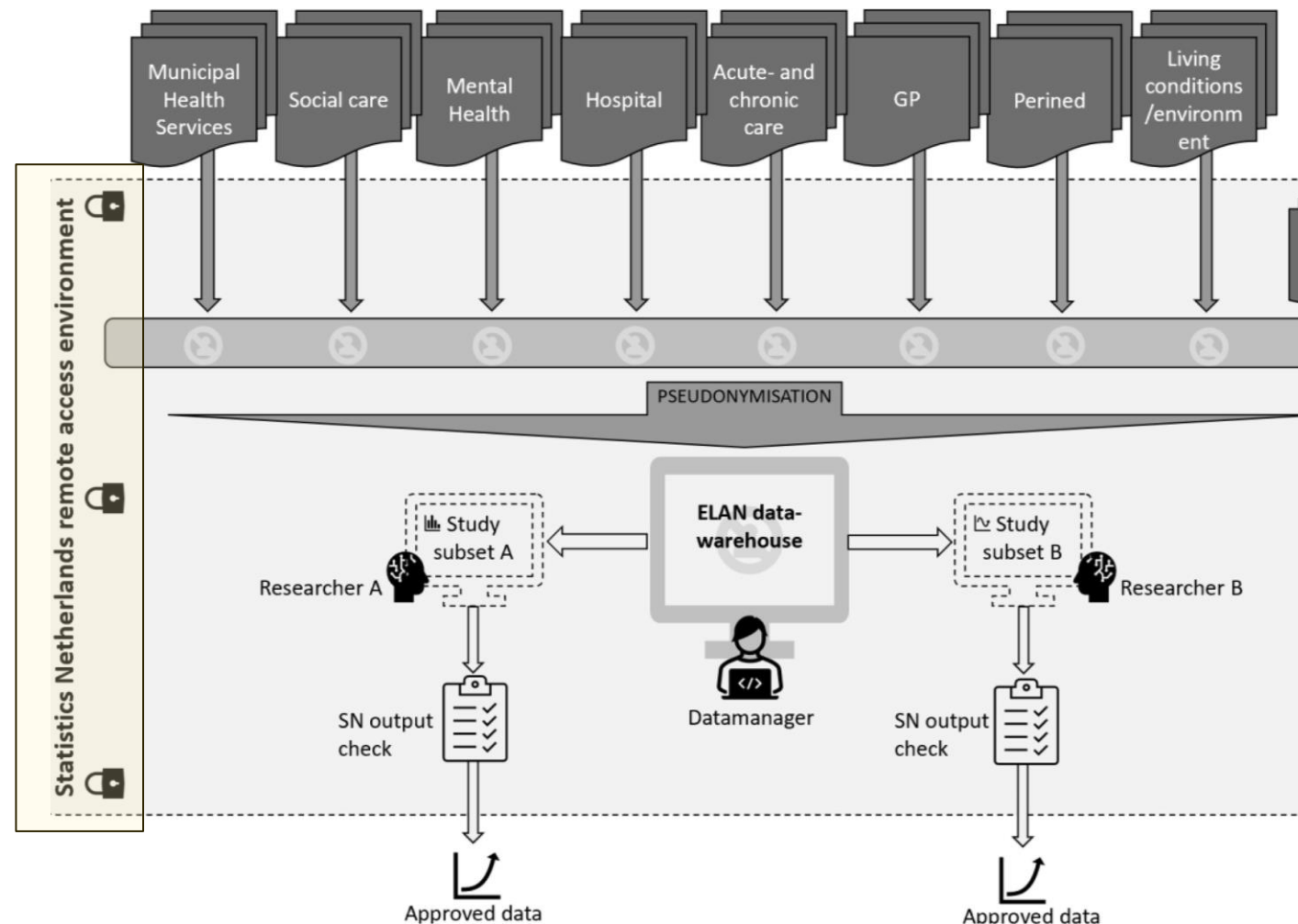
Over... ELAN

- Het “*Extramurale LUMC Academische Netwerk*” (ELAN): CBS-RA

1. trans-disciplinair
2. domein-overstijgend
3. regionale data-infrastructuur

- inclusief bijbehorende beleidsstructuur, dat gezondheidsbeleid en onderzoek ondersteunt

- Totaal aantal patiënten:
 - ~1,2 miljoen records



Health
Campus

Den
Haag



Onderwijs met synthetische data

LUMC Project Raamplan Implementatie Artsopleiding (PRIMA) 2020

NFU opdracht: *“Ontwikkelingen in technologie, waaronder informatie- en communicatietechnologie en kunstmatige intelligentie, leiden tot nieuwe mogelijkheden in preventie, diagnostiek en behandeling”*

Kans! AI-Tech met één revisie op constructieve en coherente wijze grondig integreren



Marco Spruit Karoly Szuhai

To-be: De Leiden Virtuele Patiënten & Populatie dataset (VIPP)

- Vergroot het enthousiasme van zowel studenten als docenten voor Technologie, AI & Big Data
- Door een EPD-afgeleide dataset te ontwikkelen van 'zo-goed-als-echte' patiënten en lokale bevolking

Sterktes

- Geeft impliciet een **unieke Leidse smaak** aan de GNK programma's
- Maakt **zo-goed-als-echte** data-analyses mogelijk (*i.e.* medisch betekenisvolle uitkomsten)
- Gemeenschappelijk referentiekader voor casuïstiek om **alle cursussen** te helpen verbinden

Zwaktes

- Eenmalige **ontwikkeling** is resource-intensief, vergt maatwerk (*i.e.* promovendus)
- Reeds in gebruik zijnde open datasets (*e.g.* GWAS) vereisen handmatig **koppelen** aan VIP
- Beleid tbv draagvlak behoud & medewerking van **data providers** (LUMC, GGDH)

Kansen

- Realistische data-analyses in onderwijs/onderzoek kunnen worden **getest** zonder METC
- De grote behoefte aan zo-goed-als-echte zorgdatasets zal VIP tot een nationale **gouden standaard** maken.
- De master **PHM** wil VIPP graag verweven in haar curriculum.

Dreigingen

- Vereist formeel **bewijs** dat echte patiënten nooit (opnieuw) kunnen worden geïdentificeerd
- AVG-gerelateerd **vertrouwen** garanderen, gegeven lastig uit te leggen bewijzen (XAI)
- Hoeveel inspanning is er nodig voorbij "**plausibele ontkenning**" (LaPlace-ruis)?

**Health
Campus**

**Den
Haag**

Technieken voor synthetische data

Drie benaderingen voor onderzoek en toepassingen

- A. Regel-gebaseerd → Agent-Based Modelling (ABM)
- B. Neuraal netwerk-gebaseerd → Generative Adversarial Networks (GAN)
- C. Off-the-shelf → Syntho



Ammar Faiq Jim Achterberg Marcel Haas

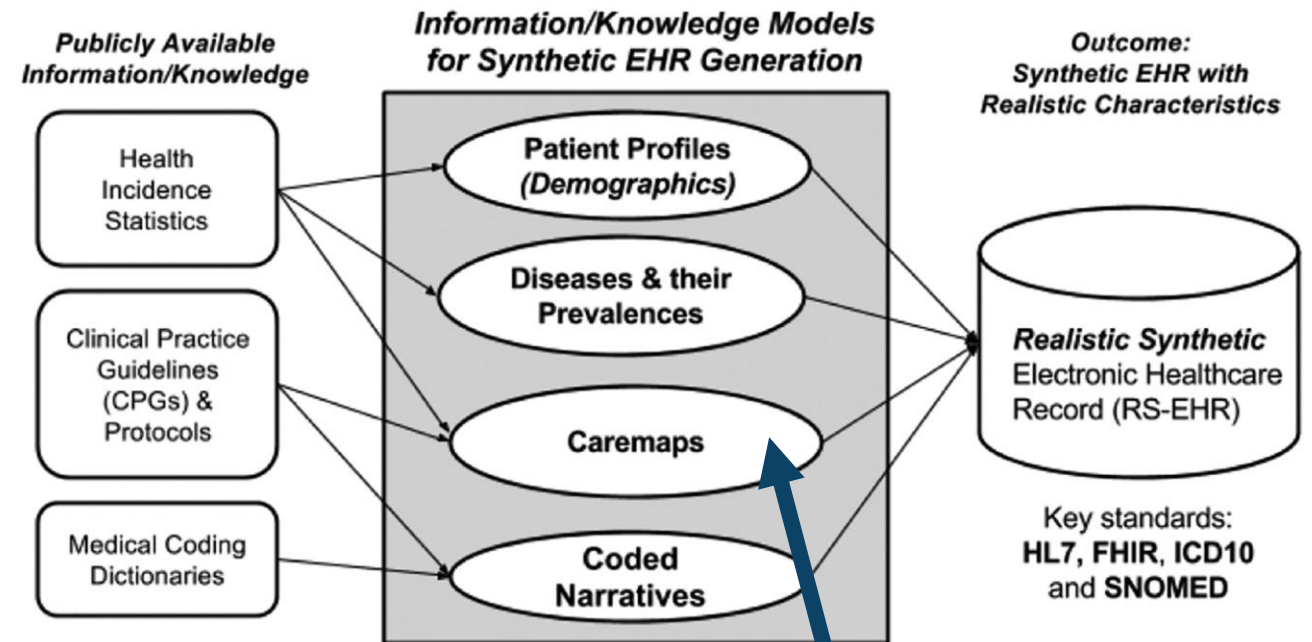
A. Agent-Based Modelling (ABM)

- *Synthea*: an open-source Synthetic Patient Population Simulator
 - Doel: “[...] simulate the lifespans of synthetic patients, modeling the 10 most frequent reasons for primary care encounters and the 10 chronic conditions with the highest morbidity **in the United States**”
 - EPD waarden gecodeerd in standaard formaten (HL7, FHIR, CCD, SNOMED)
 - <https://github.com/synthetichealth/synthea>
 - Walonoski *et al.* (2018). [reeds 246 citaties]
 - <https://doi.org/10.1093/jamia/ocx079>



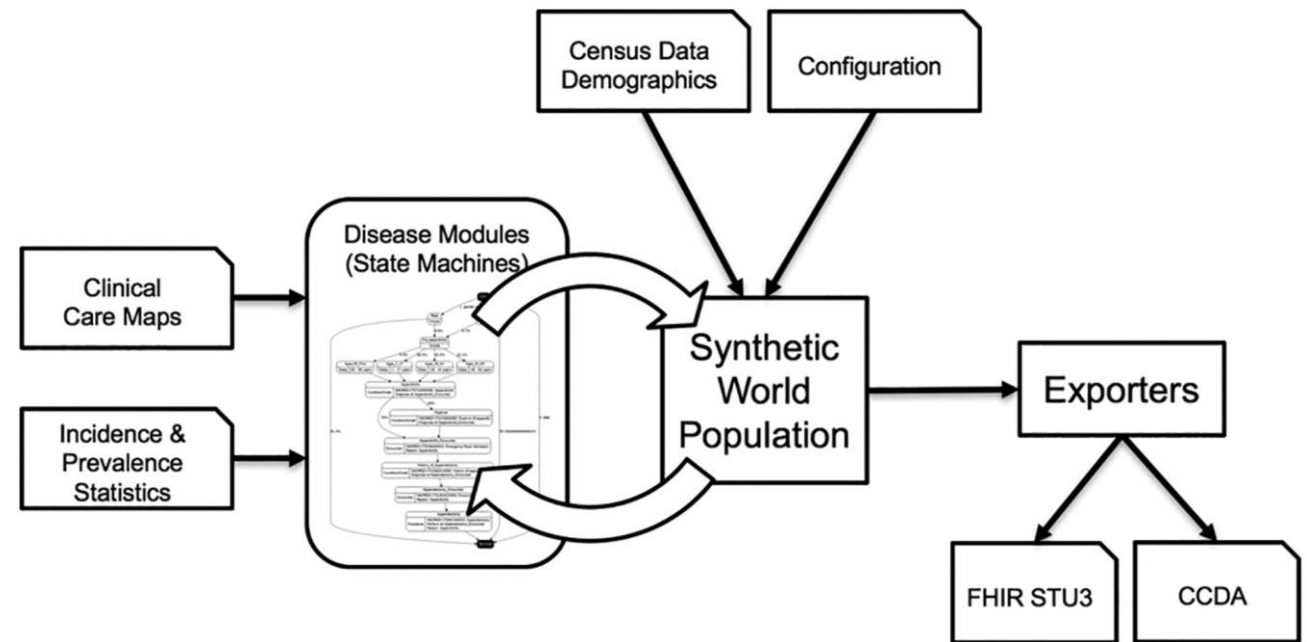
Benadering 1 “Anonymity vs Utility dilemma”

- Synthetische data (bewust) **NIET** obv originele EPD data (onwenselijk, onmogelijk), maar obv publiekelijk beschikbare data



Architectuur 1: Synthesa componenten

- Prevalentie statistieken en patiënten demografie verschillen per land maar zijn bekend (CBS!)
- Klinische zorgpaden zijn in principe mondiaal
 - mbv land-statistieken
- Ziektebeelden
 - Grotendeels mondiaal

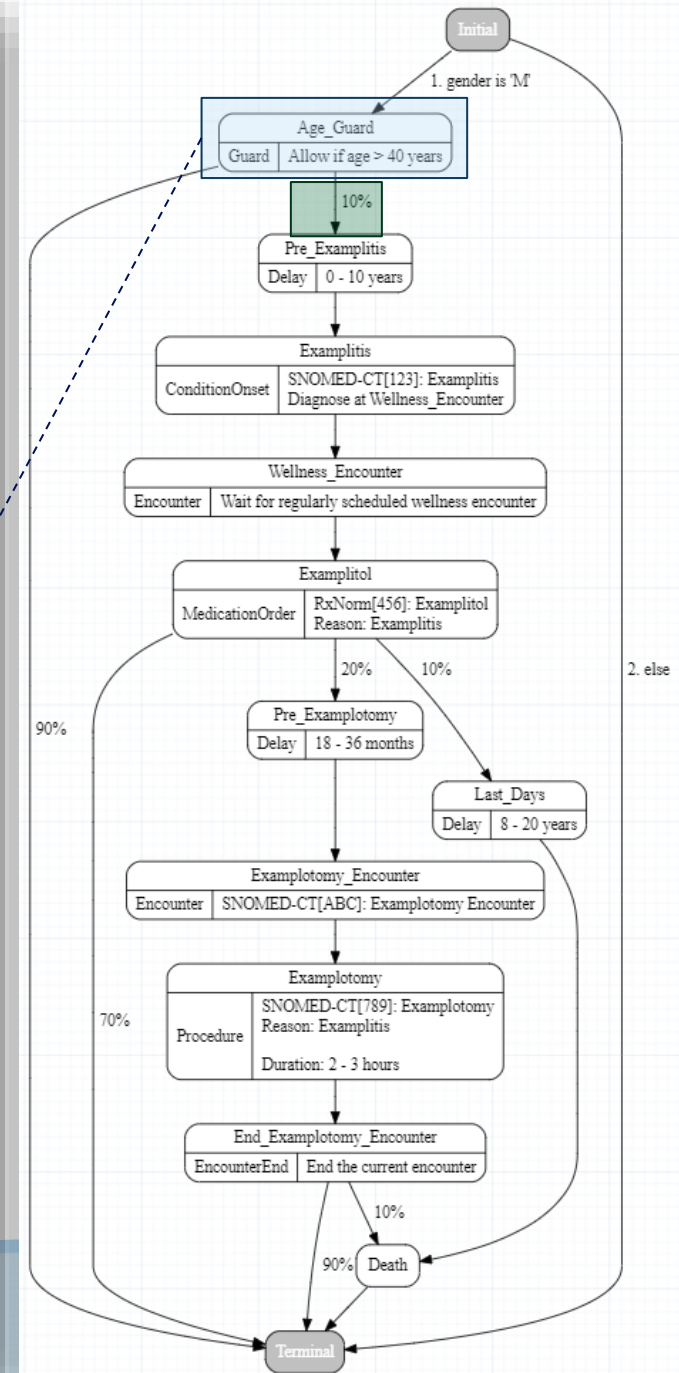


Synthea Ziektebeeld

- FSM = (“wasmachine”)
- Synthea ziektebeeld = “stochastische automaat”
oftewel een niet-deterministische eindigetoestanden-automaat
- Vanaf geboorte iedere week uitgevoerd voor iedere patiënt

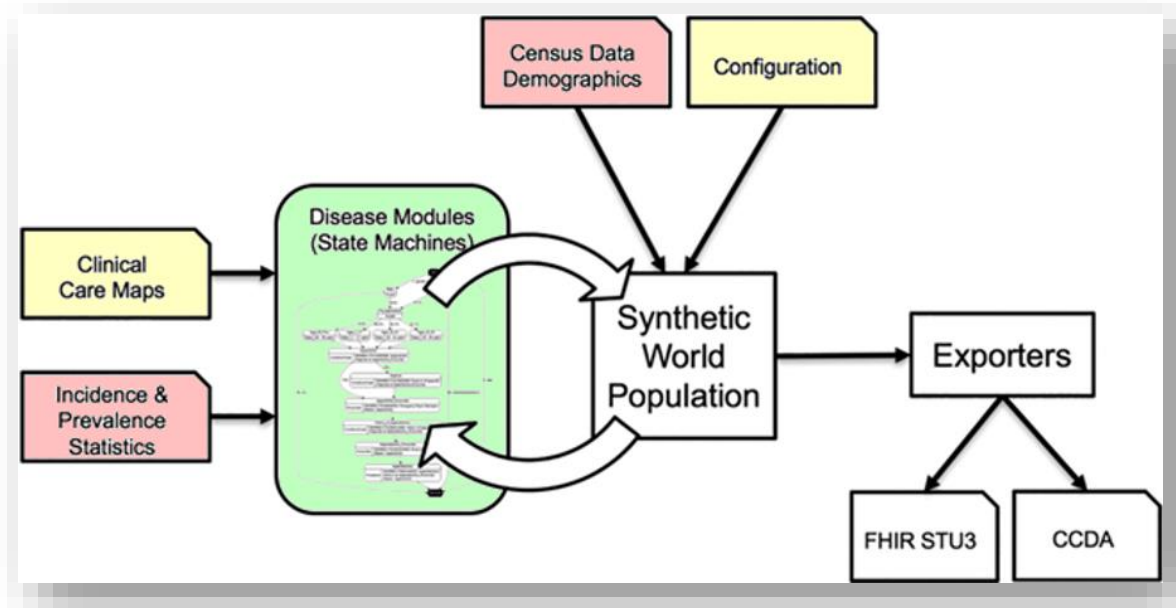
Exemplitis

```
1 {
2   "name": "Exemplitis",
3   "remarks": [
4     "Exemplitis is a painful condition
5     that affects only males. Most
6     patients ",
7     "can be cured with Exemplitol or an
8     Examplotomy but some never recover."
9   ],
10  "states": {
11    "Initial": {
12      "type": "Initial",
13      "conditional_transition": [
14        {
15          "condition": {
16            "condition_type": "Gender",
17            "gender": "M"
18          },
19          "transition": "Age_Guard"
20        },
21        {
22          "transition": "Terminal"
23        }
24      ],
25      "name": "Initial"
26    },
27    "Age_Guard": {
28      "type": "Guard",
29      "allow": {
30        "condition_type": "Age",
31        "operator": ">",
32        "quantity": 40,
33        "unit": "years"
34      },
35      "distributed_transition": [
36        {
37          "distribution": 0.1,
38          "transition": "Pre_Exemplitis"
39        },
40        {
41          "distribution": 0.9,
42          "transition": "Terminal"
43        }
44      ],
45      "name": "Age_Guard"
46    },
47    "Pre_Exemplitis": {
48      "type": "Delay",
```



Toepassing 1/2: VIPP → Synthesa voor ELAN

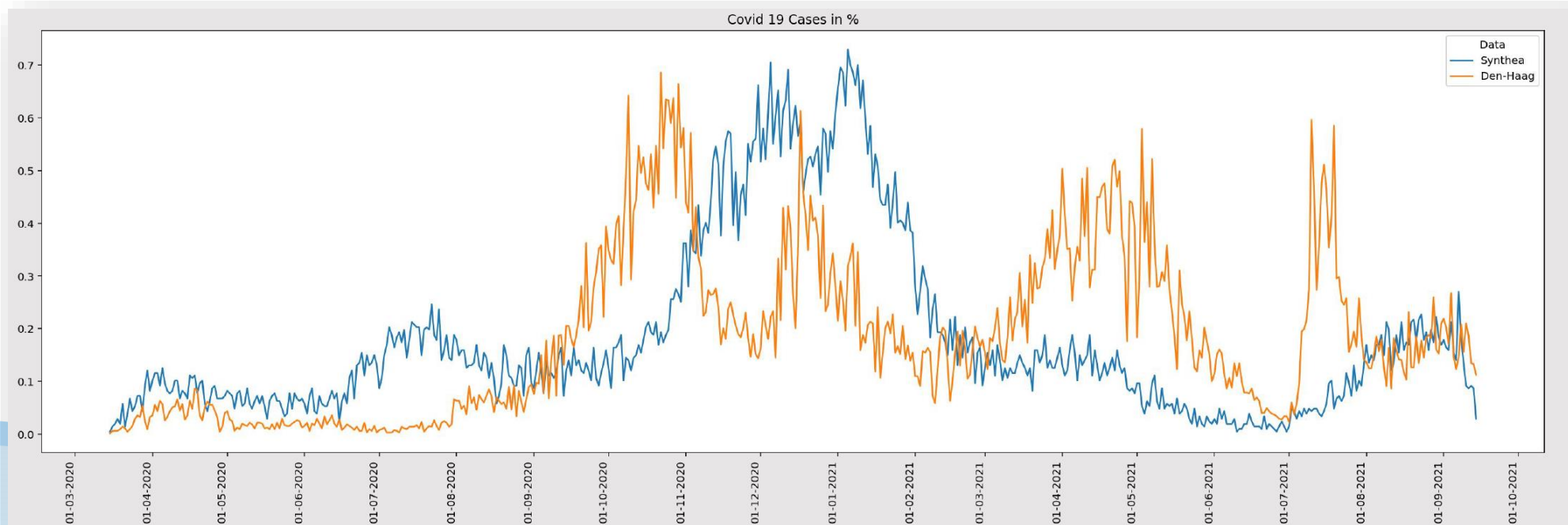
- ELAN doel: “[...] simulate the lifespans [...] in the ~~United States~~”
 - **Aanpassen:** US prevalentie statistieken en patiënten demografie naar NL
 - **Tweaken:** zorgpaden & standaarden
 - **As-is:** ziektebeelden



- Meer details bij onze VIPP poster presentatie!

Toepassing 2/2: VIPP → Synthesa voor PHM

- MSc *Fundamentals of Population Health Management* vakdoel
 - Praktijkopdracht data-analyse COVID19 voorspelmodel
 - Vereist integrale analyse van medische, sociale, regionale, etc data
- **PHM casus...** Meer details bij LUMC/VIPP poster presentatie!



Health
Campus

Den
Haag

Technieken voor synthetische data

Drie benaderingen

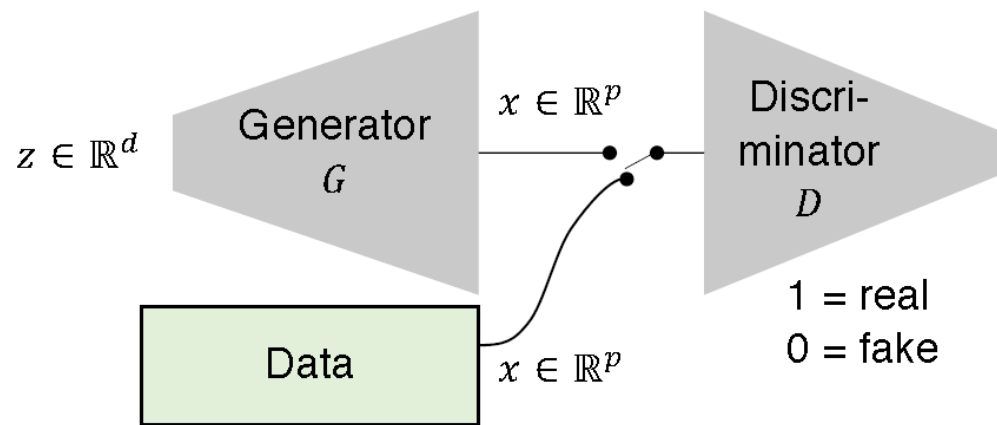
- A. Regel-gebaseerd → Agent-Based Modelling (ABM)
- B. **Neuraal netwerk-gebaseerd → Generative Adversial Networks (GAN)**
- C. Off-the-shelf → Syntho



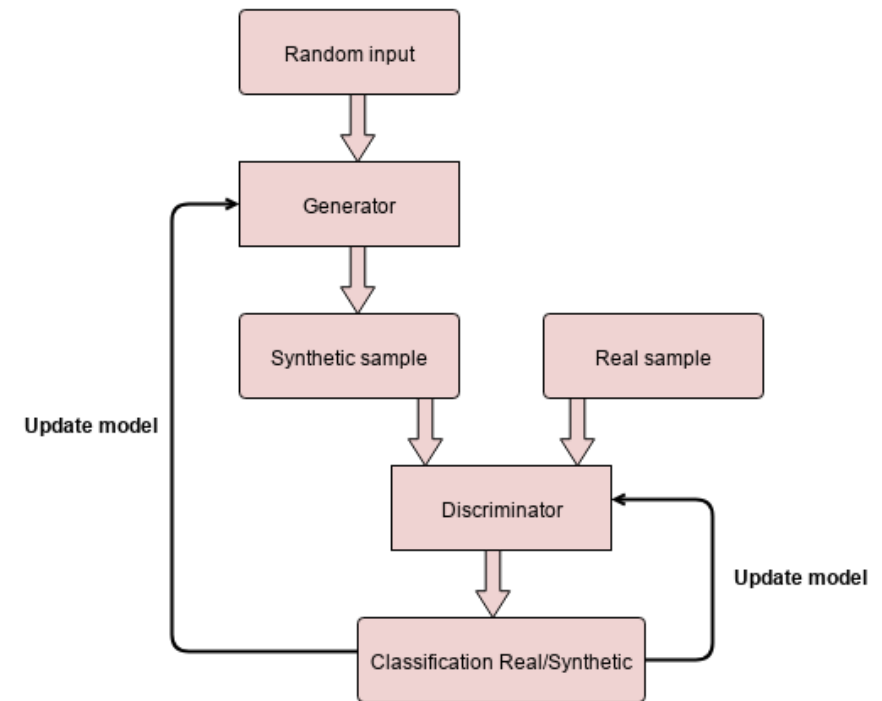
Ammar Faiq Jim Achterberg Marcel Haas

B. Generative Adversarial Networks (GAN)

- “Het doel van een generatief model is het bestuderen van
 1. een verzameling trainingsvoorbeelden, en
 2. de waarschijnlijkheidsverdeling die dit heeft voortgebracht.”



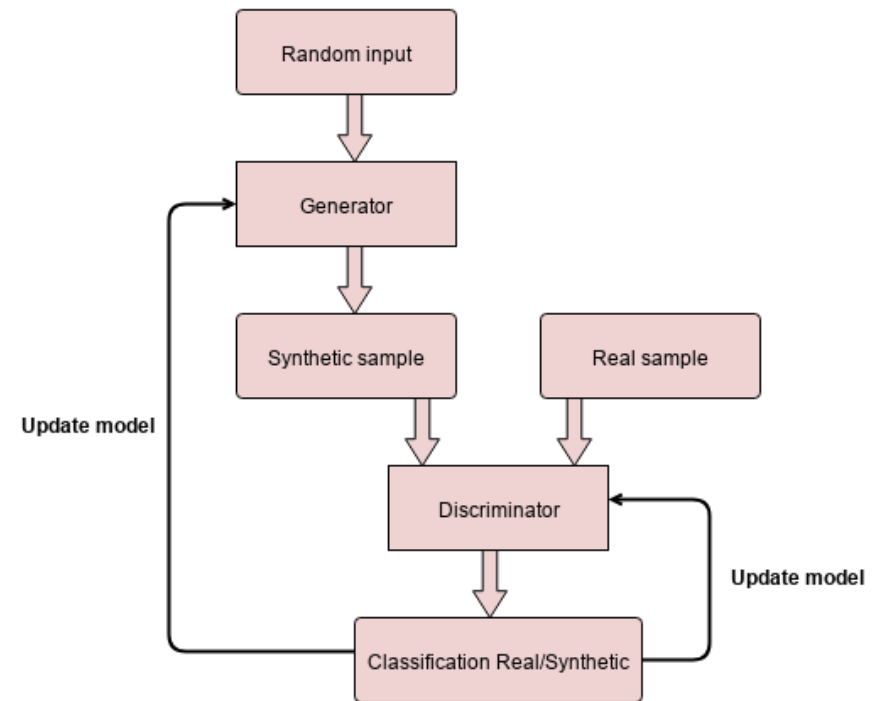
- GAN trainingsproces



B. Generative Adversarial Networks (GAN)

- **GAN:** een spel tussen twee spelers: de Generator en de Discriminator.
- **Generator:** neurale netwerk dat realistische synthetische data genereert uit willekeurige input.
 - Leert alleen van interactie met de discriminator
- **Discriminator:** neurale netwerk dat de originele data en de synthetische data probeert te onderscheiden.
 - Ontvangt zowel originele data als synthetische data geproduceerd door de generator.
- De generator is succesvol als het adequaat de verdeling van de originele data leert via interactie met de discriminator.

- GAN trainingsproces



DoppelGANger: een *Conditionele* GAN

- Ontworpen voor tijdreeksen data, en zowel continue als discrete attributen
- Adresseren tekortkomingen
 1. *Complexe* correlaties tussen tijdreeksen en de gerelateerde attributen
 2. *Lange-termijn* correlaties binnen tijdsreeksen
- <https://github.com/fjxmlzn/DoppelGANger>

Approach	Flexibility	Privacy	Fidelity
raw data	no	no	best
anonymized raw data	no	?	good
Markov model	yes	yes	bad
autoregressive model	yes	yes	bad
RNN	yes	yes	bad
GANs	yes	yes	good

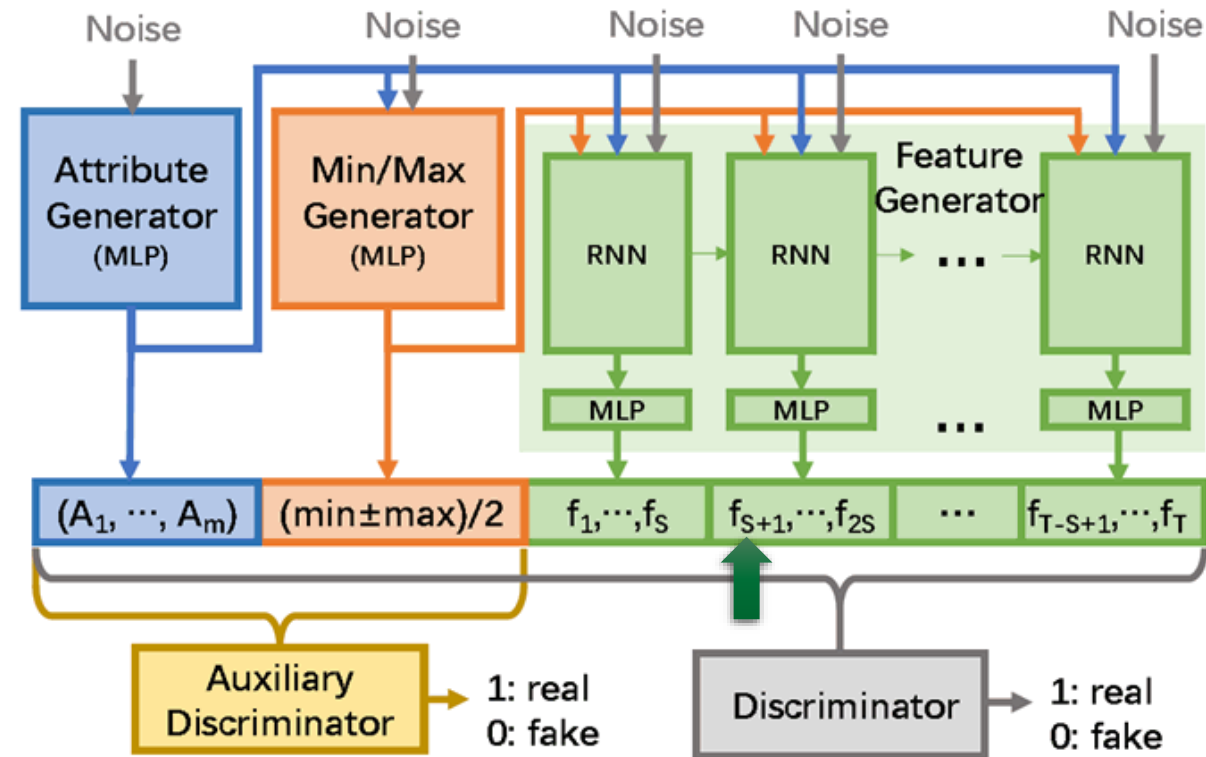
Table 1: The potential of GANs to satisfy key desirable properties of synthetic datasets.

DoppelGANger: Conditionele GAN

<klasse-specifieke distributies voor continue en discrete kenmerken>

• Top-3 kenmerken CGAN

1. Ontkoppelde normalisatie (via realistische min/max limieten 'fidelity+ if divers. → 'mode collapse')
2. Gebundelde samples generatie (om tijdsrelaties te versterken 'batch proc.')
3. Ontkoppelde attributen generatie ('fidelity+ if lange/complexere tijdsreeks')



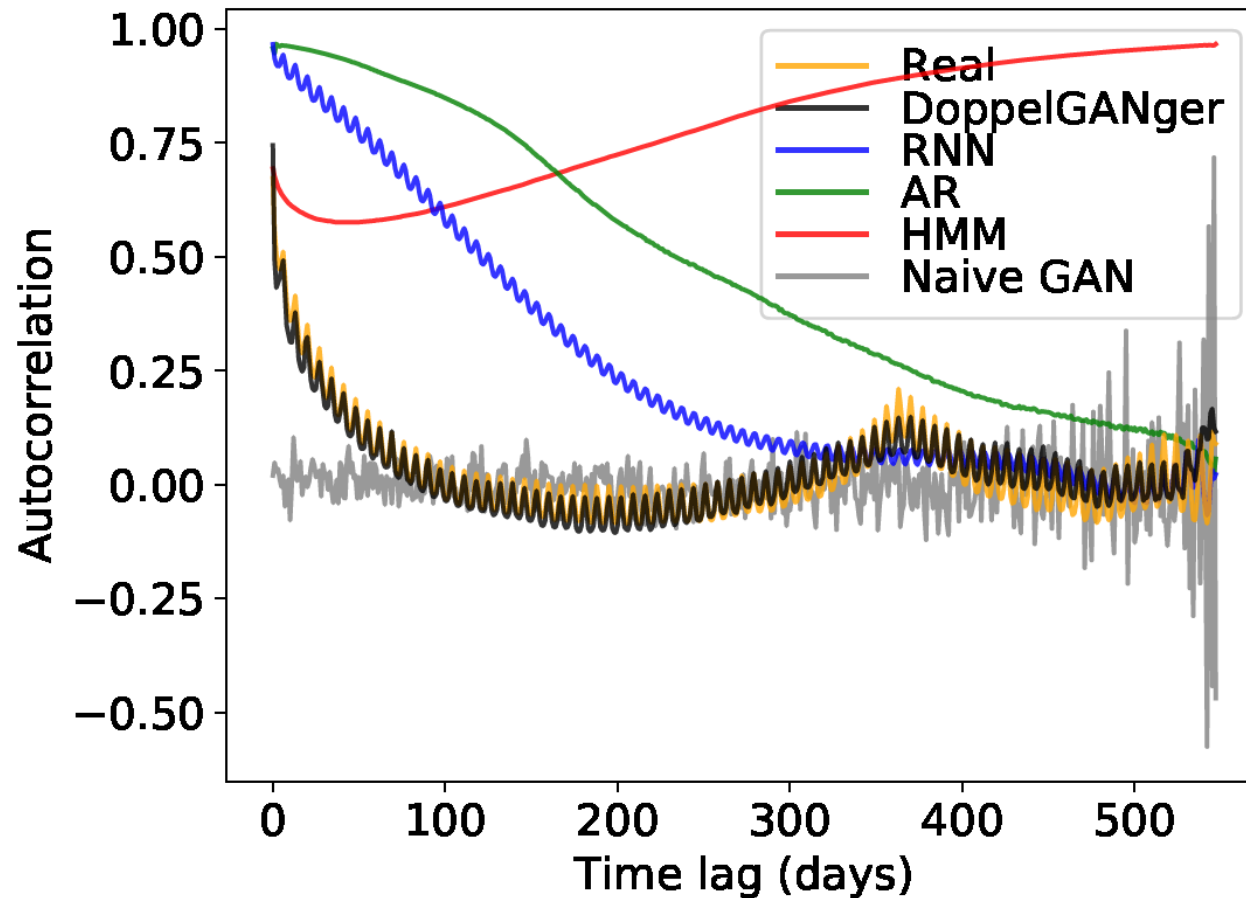


Figure 1: Autocorrelation of daily page views for Wikipedia Web Traffic dataset. DoppelGANger captures both weekly and annual correlation pattern.

Toepassing: CGAN voor ELAN (...)

- Helaas nog niet gerealiseerd... Voorbereidingen wél vergevorderd!
- Voorstudie Jim Achterberg (EUR)
 - Scheppen kader voor *evaluatie* van synthetische *medische* data: combinatie van statische en tijdsgebonden variabelen van gemengde datatypes
 - a) Uitbreiding *tSNE* algoritme tbv visualiseren tijdsreeksen van gemengde datatypes
 - b) Verbetering afstandsmeting tussen twee punten mbv *dynamic time warping & Gower*
 - c) Introductie nieuwe *two-sample goodness-of-fit test* obv classificatie-gebaseerd testen
 - d) Introductie van interpreteerbare privacy-risico maat via *attribute inference attacks*



Health
Campus

Den
Haag

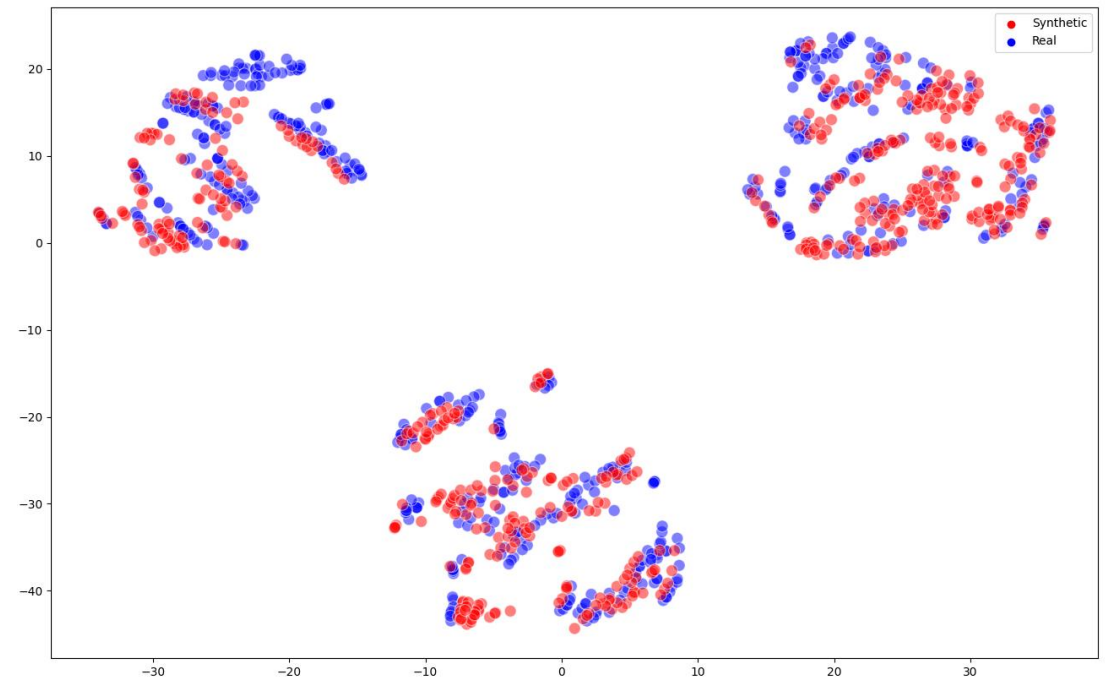
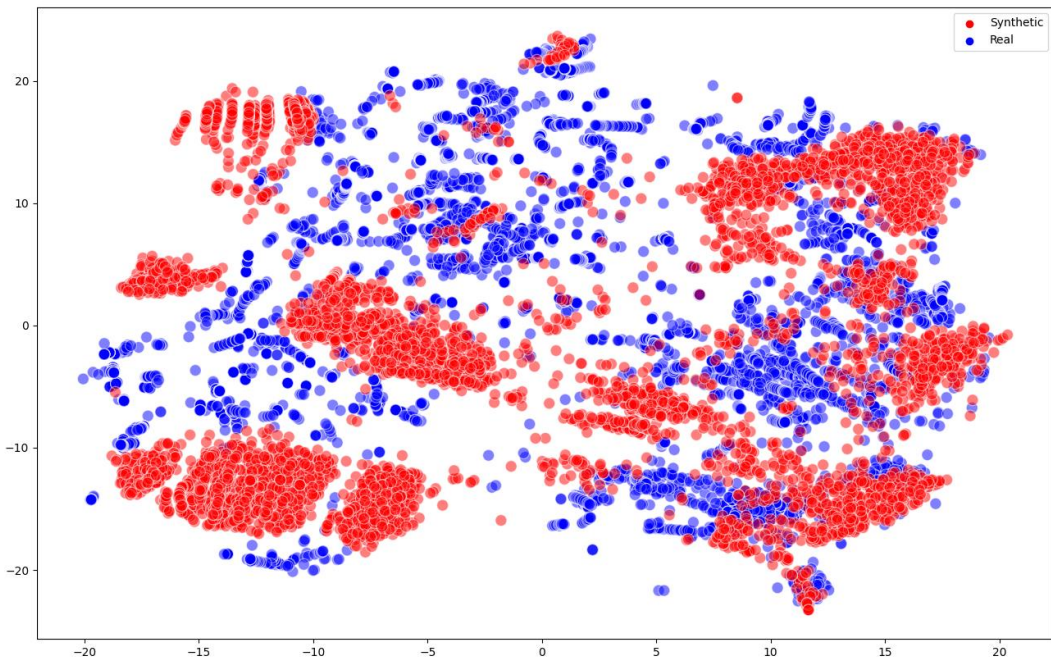
- Drie evaluatiematen
 - **Gelijkenis ‘fidelity’, Nut ‘utility’, Onthullingsrisico ‘privacy risk’**

● Synthetic
● Real

Component (a): tSNE visualisatie

<dimensiereductie via t-distributed Stochastic Neighbor Embedding>
<met behoud van locale structuur per tijdstap>

- Tijdsdimensie is impliciet, dus platgeslagen, dus tSNE toont hoe goed tijdsrelaties gegenereerd kunnen worden



Neuromuscular Monitoring (NM) dataset

- data structuur teveel 'puntenklodders'
→ te weinig data diversiteit 'mode collapse'

Cervical Dystonia (CD) dataset

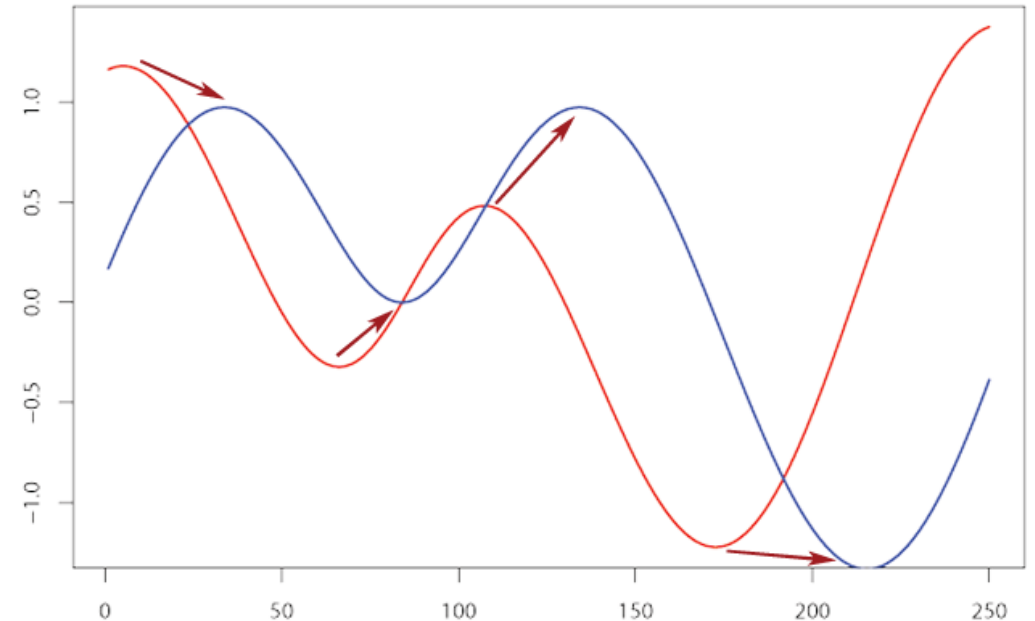
- beter behouden data structuur, maar nog niet perfect

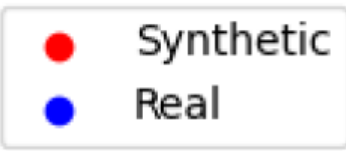
Health
Campus

Den
Haag

Component (b): Tijdsreeks afstandsmaat

- tSNE obv *Dynamic Time Warping (DTW) algoritme*:
 - Berekent afstanden tussen reeksen na uitlijning op hun gelijke vormen
 - Euclidische → Gower maat
 - (vanwege gemengde datatypes)

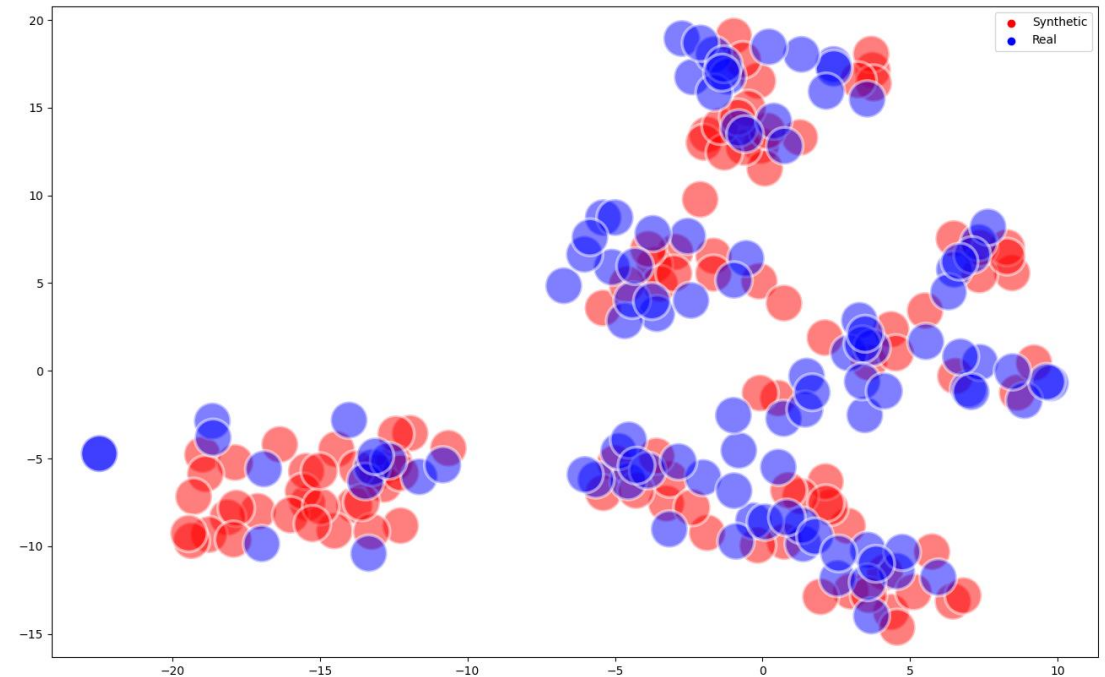
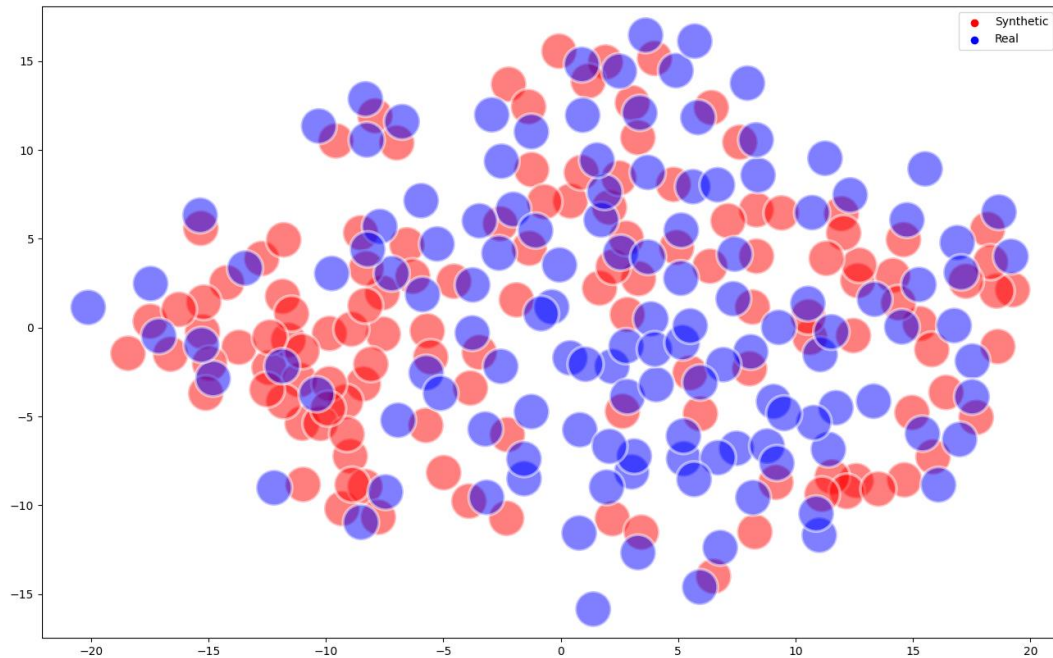




Componenten (a+b): m-tSNE visualisatie

<dimensiereductie via t-distributed Stochastic Neighbor Embedding>
 <met behoud van locale structuur per tijdsreeks (=nieuw)>

- Visualisatie van gehele tijdsreeksen ipv afzonderlijke tijdstappen, door afstanden te berekenen via DTW + Gower



Neuromuscular Monitoring (NM) dataset

- wederom enige 'puntenklodders'
- te weinig data diversiteit = 'mode collapse'



Cervical Dystonia (CD) dataset

- ook enige 'puntenklodders'

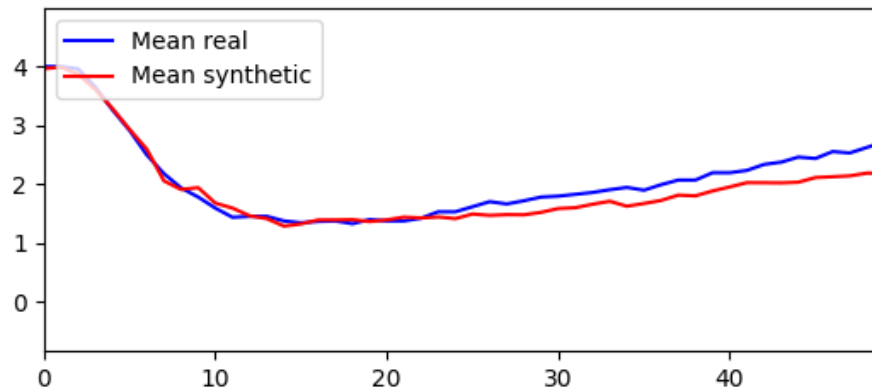
Component (c): Medische Goodness-of-Fit test

- Multi-variate *two-sample GOF* test:
 - Zijn twee distributies gelijk of niet?? Test uitvoeren!
 - Vereist wel een geschikt voorspelmodel, die samples van gemengde tijdseries en statische attributen kan beoordelen
 - Naar Friedman's (2003) intuïtie: **“als twee samples uit dezelfde distributie komen, dan kunnen ze op gelijke wijze gescoord worden door een voorspelmodel”**
 - Voorspelmodel kan hetzelfde RNN model zijn als voor evaluatie van discriminerende scores

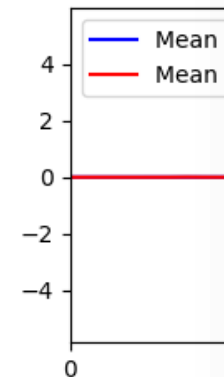
Evaluatiemaat: 'Fidelity'

- **Gelijkenis:** Gelijkenis van distributies tov originele data
 - i. *Descriptieve statistiek:* bereik, gemiddelde, standaard deviatie
 - ii. *Twee-dimensionele visualisatie:* tSNE'
 - iii. *Medical Goodness-Of-Fit (MGOF) test*

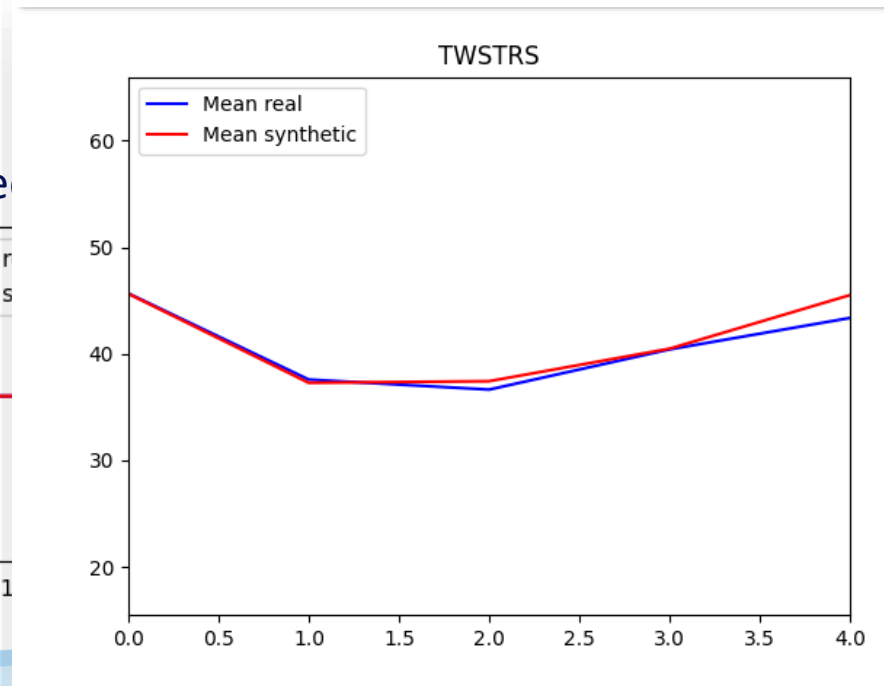
i) Aantal spiertrekkingen onder narcose



i) Hoeve



i) Gemiddelde ontwikkeling over de tijd



Evaluatiemaat: *'Utility'*

- ***Nut***: Nut bij gebruik in praktijktaken tov originele data
 - Voorspelmodel trainen op synthetisch vs originele data → vergelijken
 - *Groot nut =?* Synthetische data is bruikbaar in praktijk ipv originele data
 - i. Train Synthetic, Test Real (TSTR)
 - ii. Train Real, Test Real (TRTR)
- ***“Gouden standaard”***: via reeds gepubliceerde ELAN data-analyses
 - Niet enkel vergelijken met originele data, ook met origineel voorspelmodel

Evaluatiemaat: 'Privacy risk'

- **Onthullingsrisico:** Nut bij gebruik in praktijktaken tov originele data
 - 'Attribute Inference Attack (AIA)': Voorspelmodel trainen op synthetische data, om gevoelige attributen uit originele test-data te achterhalen

- Origineel=TRTR, GAN=TSTR
- Bij AIA accuratesse:
 - laag-laag: predict--, privacy++
 - hoog-laag: fidelity--,
 - **hoog-hoog:** predict++, privacy—
 - ~50% = willekeurige voorspelling
→ weinig risico op onthulling

Maat	Attribuut	Model	Dataset	
AIA Accuratesse	Geslacht		NM	CD
		Origineel	0.692	0.690
		GAN	0.558	0.55
Gemiddelde fout (MAPE)	Leeftijd		NM	CD
		Origineel	0.392	0.155
		GAN	0.695	0.193
	BMI		NM	
		Origineel	0.181	
		GAN	0.21	

**Health
Campus**

**Den
Haag**

Technieken voor synthetische data

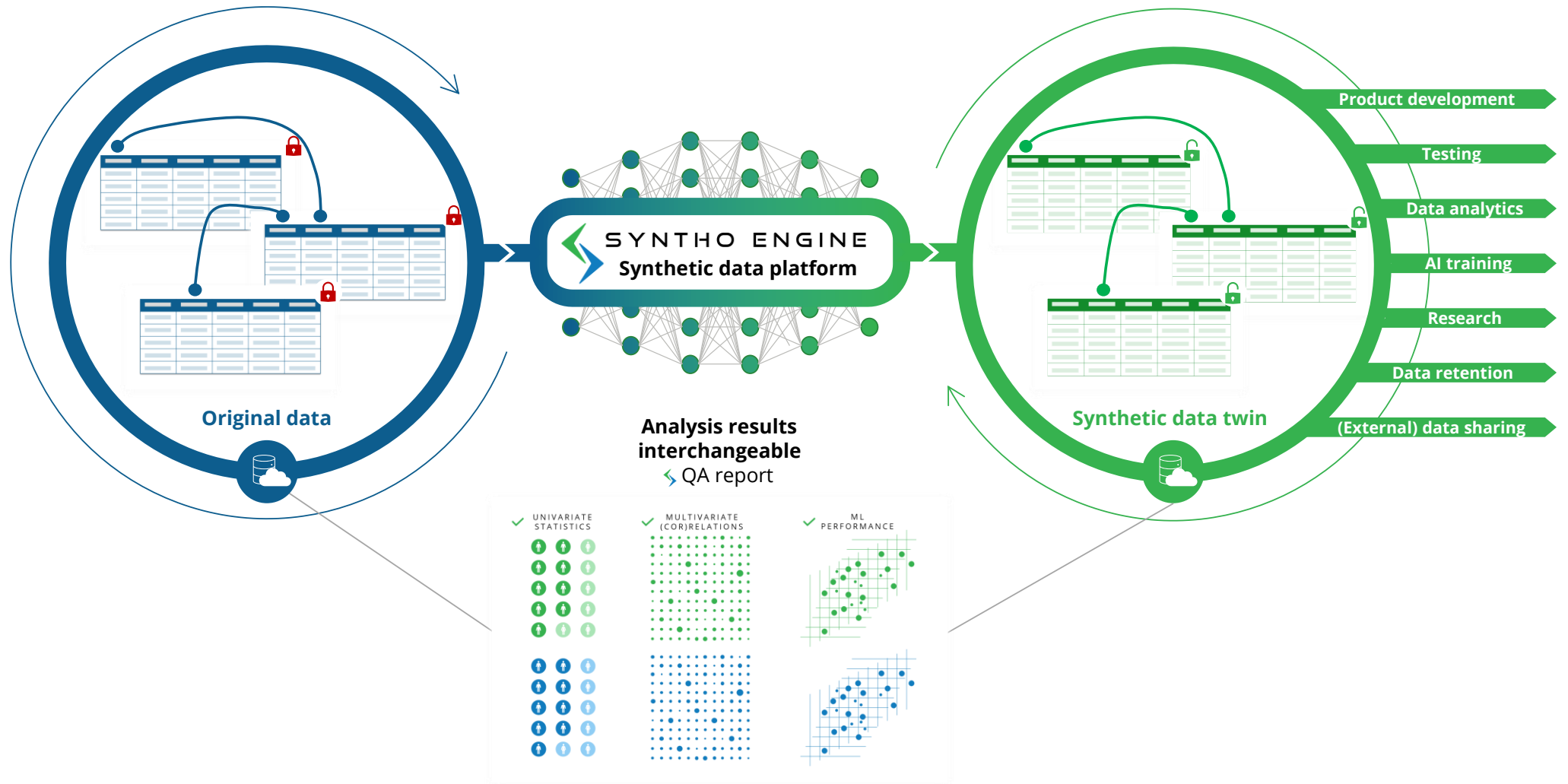
Drie benaderingen

- A. Regel-gebaseerd → Agent-Based Modelling (ABM)
- B. Neuraal netwerk-gebaseerd → Generative Adversarial Networks (GAN)
- C. Off-the-shelf → Syntho

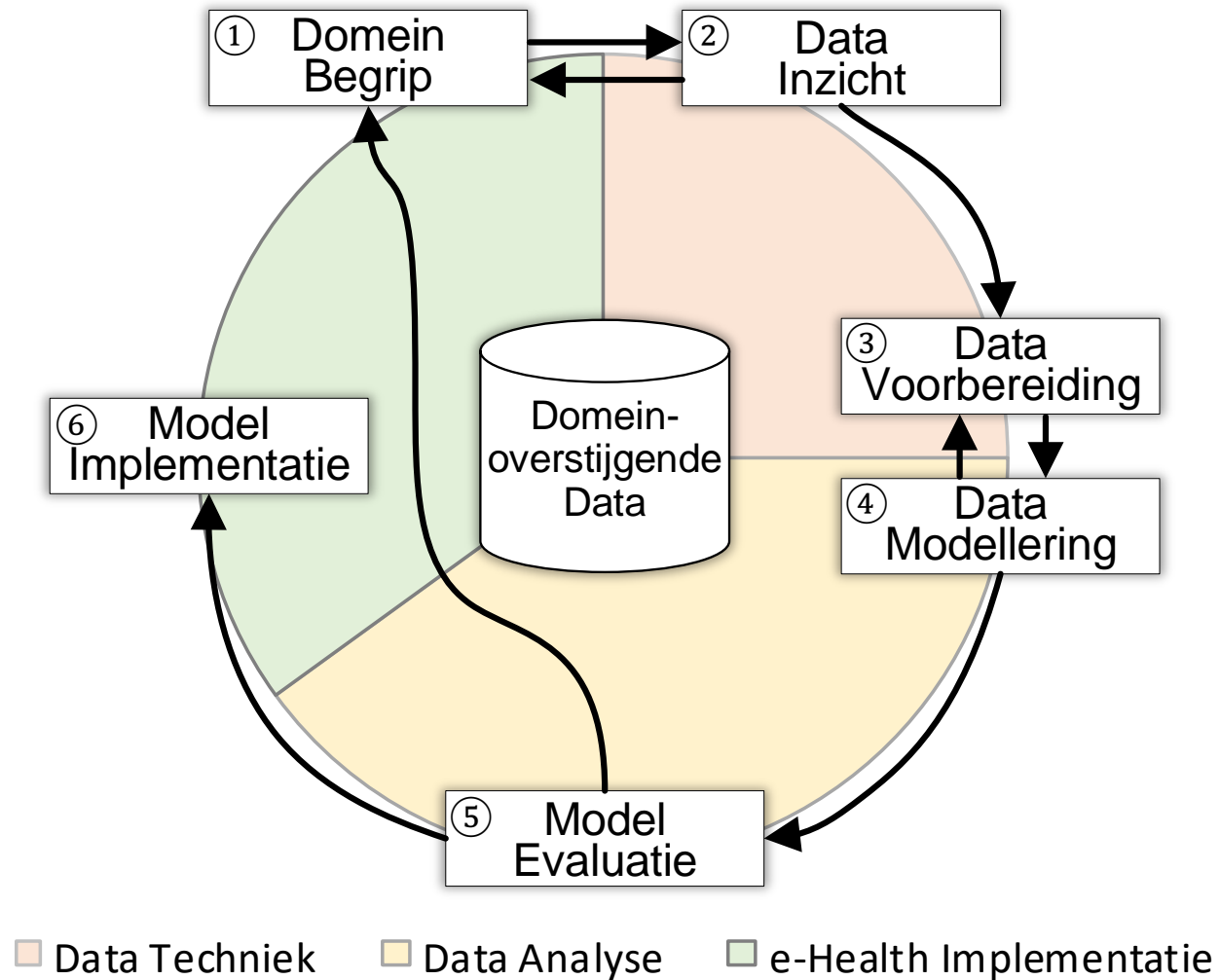


Ammar Faiq Jim Achterberg

C. Off-the-shelf (Syntho *platform*)



C. Off-the-shelf (Syntho *platform*)



**Health
Campus**

**Den
Haag**

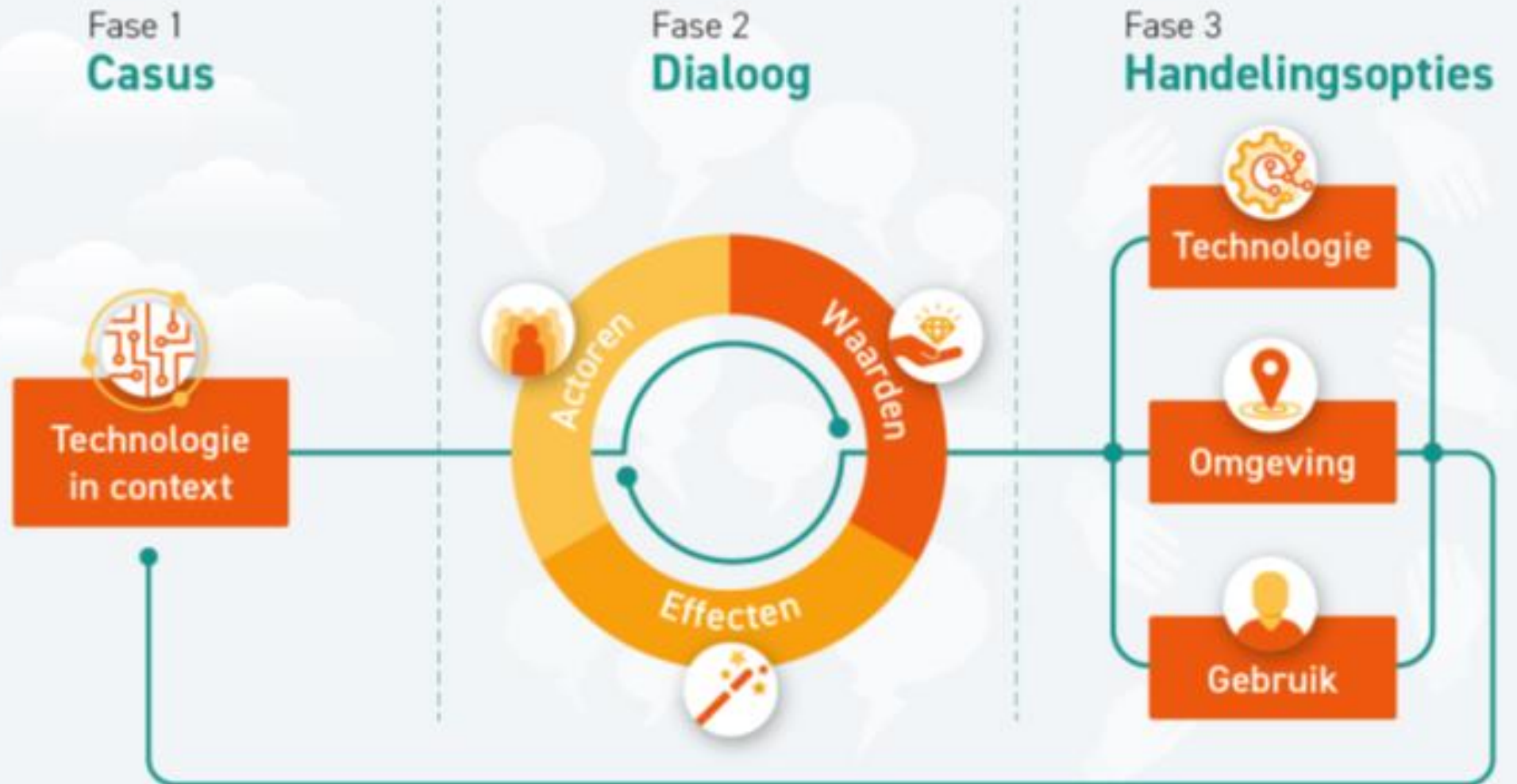
Et cetera

Effecten en de nabije toekomst



André Krom

Aanpak begeleidingsethiek



Nabije toekomst

- **CPAN**: Probabilistische Autoregressieve Netwerken (PAN)
 - als regulier neurale netwerk, maar met vertraagde invoer als tijdsmodel
- **Tekstdata** generatie: ABM+GAN = chatty-notes
 - Synthetische *FHIR Bundle* → episode extractie → ChatGPT prompt
 - Eerste GPT-2 poging met MI afstudeerder helaas voortijdig gestrand
- Samenwerking met promotieonderzoek Universiteit van **Murcia**
 - Onderzoek gaande op Madrileense EPD data (2 ziekenhuizen)
- **Nu**: projectsamenwerking met...

Health
Campus

Den
Haag

Dank! m.r.spruit@lumc.nl

“Op naar een digitale tweeling mét ELAN”

VIPP-CBS is een samenwerkingsverband tussen CBS, LUMC, Universiteit Leiden, Syntho en Health Campus Partners



SYNTHO

