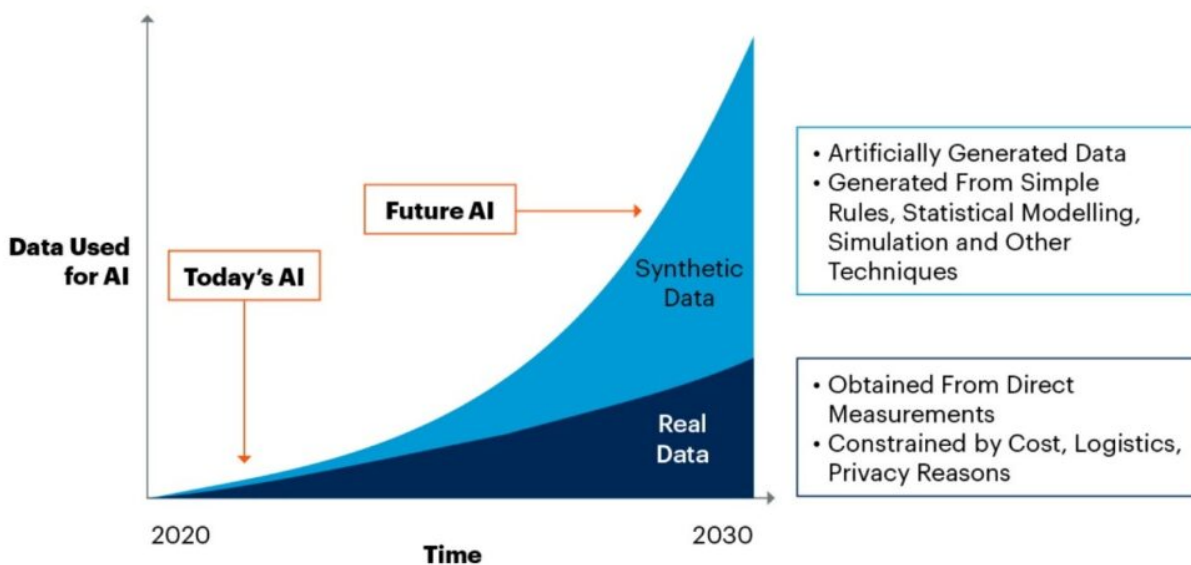


## Synthetische data voor AI

Het goed trainen van AI-modellen vereist data. Vél data. Afhankelijk van de situatie kan het lastig zijn om binnen een redelijke tijd een goed volume aan trainingsdata te verzamelen. Het verzamelen van trainingsdata is extra lastig in situaties waarbij de beslissingen van de AI-modellen effect hebben op de data die je meet. Denk bijvoorbeeld aan de situatie waarin een AI-model gebruikt wordt voor het aansturen van rioolgemalen. Op basis van de data over waterstanden besluit het model om wel of niet te pompen, wat weer effecten heeft op de waterstanden.

Vanwege die moeilijkheden wordt er steeds vaker gewerkt met 'synthetische' data. Synthetisch data zijn kunstmatig gegenereerde data die (hopelijk!) goed genoeg lijken op de praktijkdata. De verwachting is dat in de toekomst er zelfs meer gebruik gemaakt zal worden van synthetische data dan van echte data (Figuur 1).

### By 2030, Synthetic Data Will Completely Overshadow Real Data in AI Models



Source: Gartner  
750175\_C

Gartner

Figuur 1: Toename van synthetische data

## Synthetische data door middel van simulaties

Voor het genereren van synthetische data kun je allerlei tools of zelfgebouwd computermodellen gebruiken. In deze post gaan we in op een specifiek geval van het genereren van data, namelijk het gebruik van **simulaties**.

Bij simulaties bouw je een computerprogramma's dat de praktijk nabootsen. Een simulatie levert dan data zoals die in de praktijk gemeten zouden kunnen worden, en reageert op realistische wijze

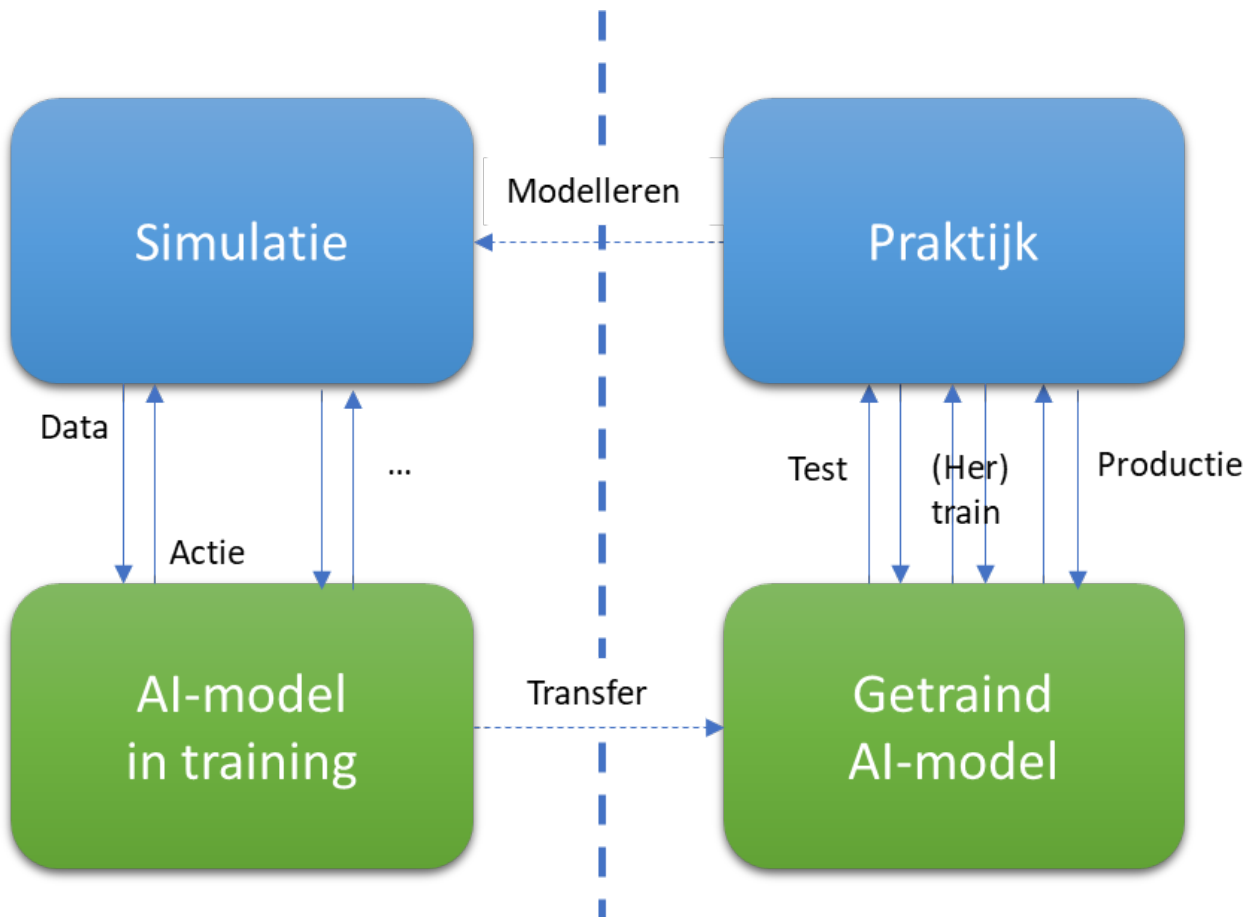
Neem als voorbeeld een simulatie van een computernetwerk in een bedrijfsomgeving waarop we een AI-model willen trainen om het netwerkverkeer te regelen. In dat geval moet de simulatie dus een patroon van netwerkverkeer leveren dat typisch is voor het soort netwerk in een dergelijke bedrijfsomgeving. Verder wil je dat de simulatie na beslissingen van het AI-model (bijvoorbeeld 'stuur deze datastroom langs die route door het netwerk') de effecten laten zien zoals die daadwerkelijk in de praktijk optreden (bijvoorbeeld dat het netwerk beter functioneert voor bepaalde gebruikers).

Het gebruik van simulatie kan veel voordelen bieden:

- Het is een goedkope en snelle manier om aan data te komen.
- Je kunt geautomatiseerd vele scenario's trainen, inspecteren en herhalen.
- Je kunt in een simulatieomgeving situaties die in de praktijk zelden voorkomen, maar wel een grote impact hebben, simuleren.
- Je kunt makkelijker gebalanceerde trainingssets maken. Sommige classificatie-algoritmen hebben namelijk een sterke voorkeur voor trainingssets waarbij alle mogelijke classes (grotweg) even vaak voorkomen.
- Je voorkomt risico's van het loslaten van onvoldoende getrainde algoritmes op de praktijk.
- En mocht je nog niet weten welk soort AI-model het meest geschikt is voor je toepassing, dan kun je in de simulatie verschillende modellen tegen elkaar uittesten.

Figuur 2 laat de typische elementen zien bij de inzet van simulatie. Op basis van praktijkdata en expertkennis van de praktijksituatie wordt een simulatiemodel gemaakt dat (hopelijk!) de relevante kenmerken van de praktijk goed kan nabootsen. Vervolgens kan het AI-model de simulatie gebruiken om te trainen. Het AI-model ontvangt data uit de simulatie en neemt op basis daarvan acties. Die acties worden weer aan de simulatie teruggeven wat indien nodig weer leidt tot aanpassingen in de toestand van het simulatiemodel.

Wanneer de training is afgerond, of in elk geval het maximale leereffect is bereikt, dan wordt het model in de praktijksituatie gebracht (transfer). Het is dan belangrijk om het model beheerst te testen en eventueel aanvullende training te doen, voordat het in productie wordt genomen. Uiteraard is het modelleren, trainen, transfereren en testen in de praktijk een iteratief en cyclisch proces.



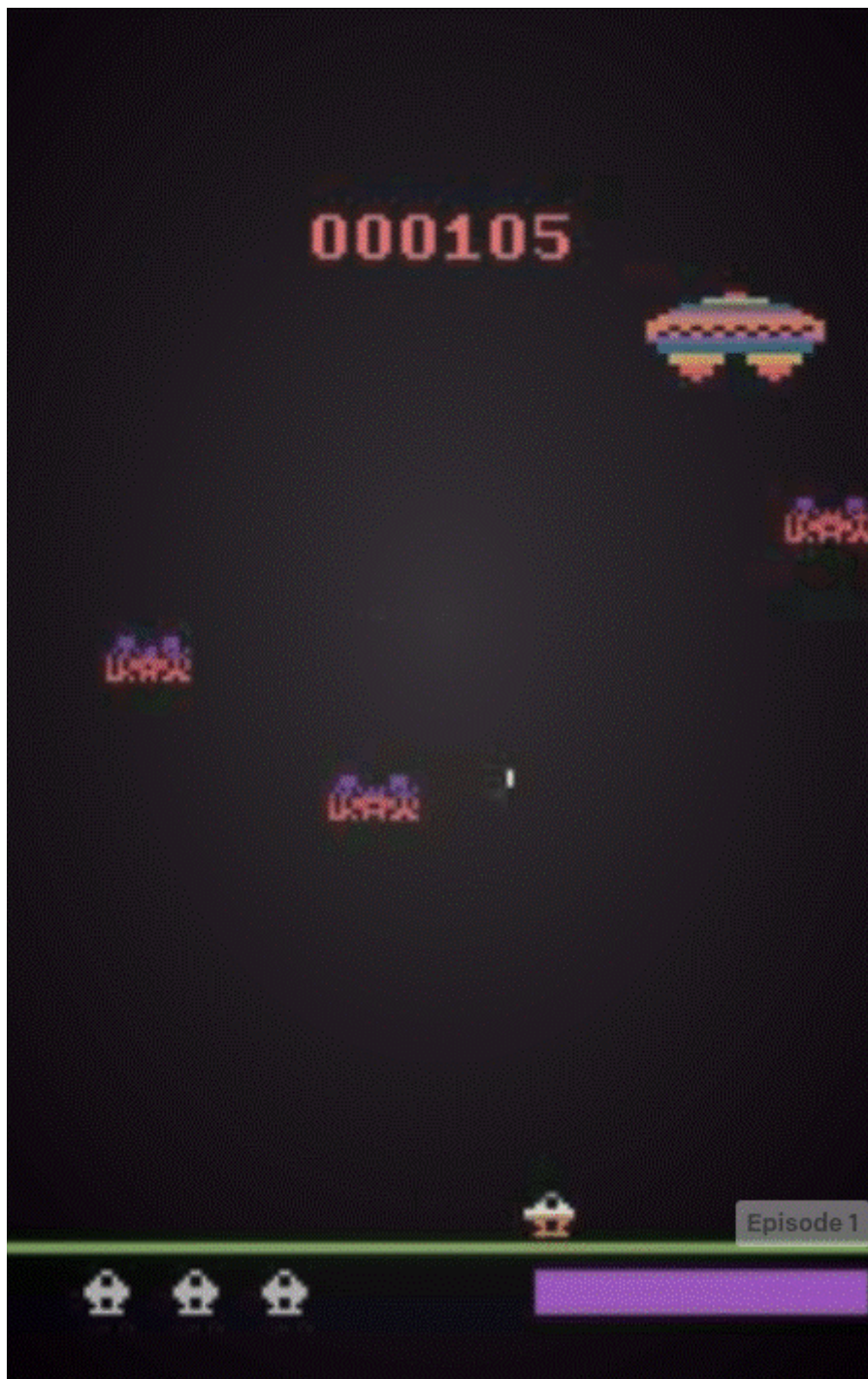
Figuur 2: Van praktijk naar simulatie en weer terug

## Wanneer is de inzet van simulatie zinvol?

Simulatie is echter niet in alle situaties even goed bruikbaar. De belangrijkste factor hierin is hoe goed het mogelijk is om het gedrag van de praktijk na te bootsen in de simulatie. We bekijken een drietal scenario's.

### 1. De praktijk laat zich exact nabootsen in de simulatieomgeving.

Dit is het ideale geval waarbij alle variabelen onder controle zijn. Een typisch voorbeeld hiervan is bij computerspellen, zoals bijvoorbeeld de klassieke Atari-games. Als de simulatieomgeving en praktijkomgeving gelijk zijn, dan kan het AI-model trainen in dezelfde omgeving als waarin de uiteindelijke prestatie wordt gemeten. De beste prestatie in de training is dan meteen ook het eindresultaat.



Figuur 3: Een Atari game zoals te vinden op OpenAI gym

- **De praktijk laat zich ‘voldoende nauwkeurig’ nabootsen in de simulatieomgeving.**

Ook in dit geval is simulatie goed bruikbaar. Na het trainen wordt het model getransfereerd naar de praktijk en kan daar goed worden ingezet.

Het is makkelijker om de praktijk voldoende nauwkeurig te benaderen wanneer beslissingen van het AI-model geen effect hebben op de praktijk. Dan is de simulatie namelijk niets meer dan het nabootsen van een stroom praktijkdata. Denk hierbij bijvoorbeeld aan het simuleren van de aandelenbeurs of weerberichten.

- **De praktijk laat zich niet volledig nabootsen in een simulatiemodel.**

In deze veelvoorkomende situatie is de simulatie geen ‘tweeling’ van de praktijk, maar eerder een ‘broertje of zusje’. Dit kan gebeuren wanneer het te bewerkelijk is of simpelweg te moeilijk is om

hebben op het proces, maar te complex of ondoorzichtig zijn om mee te nemen.

In dit geval is het AI-model na training nog niet goed genoeg om in productie te draaien. Het model zal dan verder moeten leren door te oefenen in de praktijk. Bij het ontwikkelen van een robotarm kan het AI-model bijvoorbeeld een deel van de taken leren in een simulatie-omgeving en daarna verder leren met een fysieke robotarm. De simulatie heeft dan een deel van de praktijkoefening bespaard.

## Klinkt goed, wat is de catch?

In het voorgaande hebben we gezien dat de bruikbaarheid van simulatie afhangt van hoe goed het lukt om de praktijk te vangen in het model. De catch zit hem er echter in **dat het vóóraf vaak onbekend is of het simulatiemodel voldoende nauwkeurig zal zijn**. Je weet vaak niet welke data je moet simuleren opdat het getrainde model goed zal werken in de praktijk. Heb je in het voorbeeld van het AI-model dat handelt op de aandelenbeurs genoeg aan het simuleren van vraag en aanbod? Of moet je voor realistische koerspatronen ook andere economische ontwikkelingen meenemen, of zelfs de uitslag van het EK-voetbal of het weerbericht?

Een verdere complicerende factor is dat je tegelijkertijd sleutelt aan de simulatiemodel en aan het AI-model. Je moet immers ook nog achterhalen wat een goed AI-model is. Er zijn dus meerdere variabelen tegelijk.

Vaak zie je pas bij het testen in de praktijk of de simulatieomgeving 'voldoende nauwkeurig' was. Houd daarom rekening met het aanpassen van het simulatiemodel op grond van de uitkomsten van de praktijktesten.

### Praktijktesten

Voordat je weet of het AI-model voldoende presteert in de praktijk zul je het eerst moeten testen.

Twee manieren om het AI-model beheerst te testen in de praktijk

zijn **schaduwdraaien** en **supervisie**. In beide gevallen wordt het model real-time gevoed met praktijkdata (bijvoorbeeld actuele aandelenkoersen). Bij schaduwdraaien worden de beslissingen niet daadwerkelijk uitgevoerd, maar vergeleken met de beslissingen van het bestaande systeem of mens. In het geval van supervisie worden de beslissingen wel uitgevoerd en is er een expert die de beslissingen monitort en waar nodig corrigeert.

### (Her)trainen in de praktijk

Als het simulatiemodel onvoldoende nauwkeurig is, dan kan het nodig zijn om het AI-model verder te trainen in de praktijk. Een mogelijkheid hiervoor is zogeheten 'transfer learning'. Daarbij worden de belangrijkste onderdelen uit het model hergebruikt en worden de buitenste lagen hertraind.

Een andere mogelijkheid is om de feedback van het eerdergenoemde schaduwdraaien/monitoring te verwerken in het model. En voor sommige modellen is het maar de vraag of je die ooit volledig los kunt laten op de praktijk. Bij het regelen van rioolgemaal door een AI-model kun je moeilijk tegen de bewoners van een ondergelopen wijk zeggen "jammer, maar het model heeft ervan geleerd hoor!"

Tot slot is er dan nog het concept van 'continual learning' of 'online learning'. Daar leert het model terwijl het in productie is. Dat is ook heel belangrijk als data met de tijd verandert, bijvoorbeeld als er

dat is voor een andere keer.