

Kunnen wij elke machine verslaan? Beschouwingen rondom Lucas' Argument

Albert Visser

Weinig zaken zijn fascinerender dan een kort, op het eerste gezicht eenvoudig, schijnbaar waterdicht argument met een grootse conclusie. Het mooiste voorbeeld is nog steeds Anselmus' Godsbewijs. Dat gaat ongeveer zo:

De dwaas zei in zijn hart: er is geen God. Maar zelfs de dwaas zal moeten toegeven dat hij zich iets kan voorstellen zodat niets groter is dan dat iets. Maar als dit iets niet echt, maar alleen in het verstand bestaat, kunnen we ons iets nog groters voorstellen: namelijk iets dat óók nog bestaat. We hebben nu duidelijk een tegenspraak. Dus bestaat er echt iets zodat er niets groters gedacht kan worden. Dit iets noemen we God... .

Het gekke is dat dit argument, hoe mooi ook, geen greintje overtuigingskracht heeft -voor mij niet, tenminste. Ook al zou ik geen flauw idee hebben wat er mis is, dan nog zou de conclusie mij na het horen van het argument absoluut niet waarschijnlijker voorkomen. Waarom? Omdat ik de intuïtie heb dat *zo'n soort* stelling niet bewezen kan worden met zo'n simpel 'logisch' argument. Er worden heel weinig specifieke inzichten aangaande God gebruikt, terwijl Hij -zo hij al bestond- toch een heel individueel iets zou moeten zijn. Anselmus laat God meer lijken op het neutrino, waarvan het bestaan voorspeld kon worden nog voor het experimenteel ontdekt was.

In dit artikel wil ik niet Anselmus' Argument bespreken maar dat van Lucas -J.R.Lucas om precies te zijn. Lucas' Argument is gericht tegen de dwaas die in zijn hart zei: de mens is een machine. Het argument is kort, maakt gebruik van eenvoudige, maar diepe wiskundige stellingen -de Onvolledigheidsstellingen van Gödel- en lijkt op weinig veronderstellingen te berusten. Lucas laat zien, of pretendeert te laten zien, dat er ten minste één ding is waarin wij elke machine de baas zijn: het Gödeliseren. Mijn eerste reactie op Lucas' Argument is net als bij Anselmus: ongeloof dat zo'n soort argument zo'n soort conclusie zou kunnen bewijzen. Er wordt een wiskundige stelling, die betrekking heeft op precies omschreven objecten, toegepast op een wijde, verwarrende, slechts vaag omliggende problematiek. Het toepassen van wiskunde in de techniek is al een verfijnde kunst, die de meeste mensen alleen met veel oefening onder de knie krijgen; het toepassen van wiskunde in de filosofie is nog veel moeilijker... .

Zoals we verderop zullen zien neemt nadere beschouwing van Lucas' Argument deze eerste indruk niet weg. Integendeel. Waarom is het dan toch de moeite waard om hier over dit argument te lezen? Wat maakt het tot meer dan een historische curiositeit? Ik denk dat het inzicht dat zulke argumenten niet werken zelf belangrijk is. Om dit uit te leggen moet ik eerst iets zeggen over de Onvolledigheidsstellingen van Gödel. Het verschijnen van de Gödelstellingen in 1930 was een culturele gebeurtenis van de eerste orde. De stellingen toonden aan dat Hilbert's formalistische programma voor het funderen van de wiskunde niet uitvoerbaar was. Ze gaven een impuls aan de ontwikkeling van begrippen als *mechanische rekenprocedure (algorithme)* en *formeel systeem*. Ook het esthetisch oogpunt is hier belangrijk: de bewijzen zijn wondermooi. Door combinatie van een aantal in wezen simpele ideeën worden diepe resultaten bereikt. Er bestond wijd en zijd de indruk dat de Gödelstellingen veel verdere en diepere implicaties hadden. Ze zouden essentiële beperkingen van machines, van het wiskundig denken, van de menselijke geest aantonen.⁽¹⁾ Deze drie beperkingen zijn heel verschillend en staan zelfs min of meer tegenover elkaar. In welk soort beperking men geloofde was afhankelijk van bepaalde vooropgezette ideeën. Tegenwoordig lijken dit soort 'implicaties' van de Gödelstellingen minder plausibel dan

voorheen. Waarom? In de eerste plaats zijn zulke indrukken nu eenmaal onderhevig aan culturele slijtage. In de tweede plaats weten wij nu veel meer en dat werkt altijd ontvullend: de Gödelstellingen zijn voor ons gewoon wiskundige stellingen geworden met allerlei generalisaties, beperkende clausules enzovoorts. In de derde plaats hebben filosofen zoals Lucas getracht deze 'implicaties' te beargumenteren: wij kunnen nu in detail zien waarom deze argumentatie niet werkt. Lucas heeft het vage ongearticuleerde idee dat de Gödelstellingen laten zien dat de menselijke geest niet machinaal is, concreet gestalte gegeven in zijn Argument. Het is belangrijk de onjuistheid van Lucas' Argument in te zien, omdat het analyseren van dat Argument het ons mogelijk maakt afstand te nemen van het idee dat de Gödelstellingen een belangrijk verschil tussen mensen en machines aantonen. Hierdoor krijgen we een beter zicht, op zowel de Gödelstellingen als de mens-machine problematiek.^(2,3)

De verdere opzet van dit artikel is als volgt: §1 is een korte introductie in de mens-machine problematiek. In §2 verschaf ik het benodigde achtergrond materiaal aangaande de Gödelstellingen. §3 bevat tenslotte een expositie en bespreking van Lucas' Argument.

§1 Een klein filosofisch bestiarium van vragen en doctrines rondom de mens-machine problematiek

Lucas' Argument richt zich tegen een mechanistische visie op de mens. In deze sectie wil ik in het kort beschrijven wat het Mechanisme behelst en ook wat het niet behelst. In het bijzonder wil ik duidelijk maken dat het Mechanisme een andere doctrine is dan het Fysicalisme. Het inzien van dit verschil is niet alleen van belang op 'hygiënische' gronden, maar ook om de volgende reden: *het Fysicalisme is een metafysische doctrine die voor veel beoefenaars van de moderne wetenschap prima facie plausibel is. De naieve verwarring tussen mechanistische en fysicalistische thesen maakt dat antimechanistische claims al bij voorbaat iets verdachts krijgen. Ze lijken zich immers te richten tegen iets dat al als evident geaccepteerd is: het Fysicalisme.* Lucas' Argument richt zich duidelijk tegen het Mechanisme, of het ook antifysicalistische implicaties heeft is mij volstrekt onduidelijk.

1.1 Het Mechanisme

Een rechtgeaarde Mechanist gelooft in een representatieve greep uit de volgende groep van thesen (en misschien ook nog wel in enkele nauw verwante thesen).

1.1.1 *Mensen zijn machines.*

Wel, wat zijn machines? Het is moeilijk je een definitie voor te stellen die niet zo begint: "Machines zijn artefacten die ... ". Mensen zijn natuurlijk geen artefacten, ergo mensen zijn geen machines. Dit laatste is natuurlijk wel een heel flauw argument. Het gaat de Mechanist niet om het ontstaans aspect maar om bepaalde substantiële overeenkomsten tussen mensen en machines. We kunnen alleen maar concluderen dat 1.1.1 in deze ongenueanceerde vorm onacceptabel is.⁽⁴⁾ Eén verbeteringsstrategie is het problematische woord 'zijn' te vervangen, bijvoorbeeld:

1.1.2 *Mensen lijken op machines.*

Een iets beter ogende variant van 1.1.2 is:

1.1.3 *Machinale modellen zijn adequaat voor de menselijke psychologie.*

In de literatuur beperkt men zich vaak tot specifiekere thesen die betrekking hebben op één bepaald menselijk vermogen, zoals:

1.1.4 *Het menselijk denken is mechanisch.* ⁽⁵⁾

Lucas' Argument richt zich het meest direct tegen deze laatste these. Het spreekt vanzelf dat bij een goede mechanistische theorie een uitleg hoort van wat 'mechanisch' nu eigenlijk is, met andere woorden: van wat nu precies de relevante overeenkomsten tussen menselijk denken en (zekere) machinale processen zijn. Bijvoorbeeld het tegenwoordig zo populaire Functionalisme is niet slechts een theorie over de menselijke geest, maar ook over de aard van machinale toestanden en processen; er is sprake van discrete toestanden, die functioneel van aard zijn, en van stapsgewijs verlopende processen.⁽⁶⁾

Er zijn natuurlijk ook mechanistische thesen waarin in de eerste plaats iets beweerd wordt over machines en niet over mensen:

1.1.5 *Denkende machines zijn mogelijk.*

1.1.6 *Het is denkbaar dat een machine tevens een persoon is.*

1.1.7 *Machines kunnen bewustzijn hebben.*

1.1.8 *Machines kunnen een geest hebben.*

Een voordeel van de machine-gecentreerde thesen is dat erin duidelijker de mogelijkheid wordt opengelaten dat eigenschappen als bewustzijn aanwezig zijn, zonder dat er sprake is van een substantiële gelijkheid met mensen. Ook Martianen zijn wellicht bewuste, denkende personen, maar ze hoeven geen goede modellen te zijn van de menselijke psychologie.

1.2 **Het Fysicalisme**

Fysicalisme kent vele varianten. Deze hebben alle gemeen dat er een niveau van werkelijkheid -het 'fysische'- geponeerd wordt met een speciale status. In reductionistische varianten wordt deze speciale status als volgt uitgelegd: alle andere -wellicht: schijnbaar andere- niveaus van werkelijkheid kunnen gereduceerd worden tot het fysische niveau. Wat 'reductie' hier inhoudt kan nog op vele manieren ingevuld worden. Laat ik om de discussie niet te veel te compliceren een liberale, niet reductionistische vorm van Fysicalisme schetsen. De theorie ziet er zo uit:

- a) Er is een niveau van werkelijkheid, het 'fysische', dat causaal gesloten is. Dat wil zoiets zeggen als dat gebeurtenissen op dat niveau gebeurtenissen op andere niveaus niet nodig hebben om plaats te vinden.
- b) Het fysische is datgene wat bestudeerd wordt in de fysica.
- c) Fysische objecten zijn bijvoorbeeld elektronen, electro-magnetische velden en zwarte gaten. Van alles is echter niet-fysisch.
- d) Op het fysische niveau komen zaken als functie en betekenis niet voor. Deze stoel waarop ik zit, is geen puur fysisch object. De stoel heeft een functie: hij is *om op te zitten*. Als hij die functie niet had, zou het geen stoel zijn. Fysische objecten hebben in die zin geen functie. Een stoel kan kapot gaan; fysische objecten gaan niet kapot, ze houden hoogstens op te bestaan. Natuurlijk heeft de stoel wel een fysische drager.⁽⁷⁾
- e) Veroorzaking wordt gemedieerd door het fysisch niveau. Bijvoorbeeld de gedachten van voorzitter Mao veroorzaakten (mede) de Culturele Revolutie. Maar noch de gedachten van Mao, noch de Culturele Revolutie zijn fysische objecten. De bewering is dat deze

veroorzaking niet kon plaats vinden zonder fysische belichaming van die gedachten, op papier, in geluidsgolven, in electro-magnetische golven.⁽⁸⁾

Thesen (a) en (e) beschrijven de speciale status van het fysische niveau: het is het *enige causaal zelfstandige*. De hier beschreven versie van Fysicalisme is goed te verenigen met het Functionalisme, dat functionele niet-fysische toestanden postuleert. Dit Fysicalisme richt zich wel tegen de opvatting (à la parapsychologie) dat er een andere werkelijkheid zou zijn met een andere oorzakelijkheid.

1.3 Fysicalisme impliceert Mechanisme niet

Niets in de hierboven besproken liberale variant van het Fysicalisme sluit uit dat machines en organismen een heel verschillend soort dingen zijn. Geen van beiden hoeven objecten op het fysisch niveau te zijn. Misschien hoort het wel essentieel bij machines dat ze een eindig aantal discrete functionele toestanden hebben, terwijl het begrip 'toestand van een organisme' geen duidelijke betekenis heeft.

Maar als we een mens nu eens molecuul voor molecuul nabouwen? In de eerste plaats is dat zo'n ver verwijderde mogelijkheid dat we, denk ik, niet erg goed begrijpen wat 'molecuul voor molecuul' nabouwen nu eigenlijk betekent. In de tweede plaats is het niet echt duidelijk dat we het resultaat van dat bouwen zouden moeten beschrijven als 'machine'. Ten derde: wie weet krijgen we wel zowel een machine als een mens, zonder dat die twee identiek zijn.

1.4. Mechanisme impliceert Fysicalisme niet

Het lijkt me bijvoorbeeld dat het Functionalisme als theorie van het mentale niet getrouwd hoeft te zijn met één of andere versie van Fysicalisme. (Het Functionalisme ontleent natuurlijk wel veel van haar aantrekkelijkheid aan haar *compatibiliteit* met het Fysicalisme.) Het begrip 'functionele toestand' is nog steeds zinnig ook als er twee causaal zelfstandige niveaus van werkelijkheid zijn (aannemende dat zoiets überhaupt denkbaar is). Waarom zouden organismen en machines niet aan beide niveaus deel kunnen hebben?

Lucas bestrijdt met zijn Argument het Mechanisme. Of Lucas' beweringen ook anti-fysicalistische implicaties hebben is mij volstrekt niet duidelijk.

§2 Machines, formele systemen en Gödelzinnen⁽⁹⁾

Om duidelijk te maken waarvoor de uitleg in deze § nodig is, geef ik eerst een schematische versie van Lucas' Argument. Het is niet de bedoeling dat de lezer nu al de verschillende vaktermen die in het argument een rol spelen begrijpt.

Lucas beweert dat wij mensen substantieel van *machines* verschillen. Er is iets waarin wij elke machine de baas zijn: *het Gödeliseren*. Het argument gaat zo:

- (i) Elke machine is een instantiatie (=belichaming) van een *formeel systeem*.
- (ii) Dus gegeven een machine die *consistent* is en in staat is om *eenvoudige rekenkunde* te doen, is er -volgens de stellingen van Gödel- een *ware rekenkundige zin* die de machine niet zal produceren, namelijk de *Gödelzin* van het systeem dat hoort bij de machine.
- (iii) Wij kunnen echter de waarheid van die Gödelzin inzien.

(iv) Ergo machines vormen geen adequaat model van de menselijke geest.

De opzet van deze § is als volgt: eerst leg ik uit wat we in deze context onder machines moeten verstaan (2.1) en wat onder formele systemen (2.2). En passant sta ik stil bij de *dubbele* relatie tussen machines en formele systemen: een machine kan gezien als een instantiatie van een formeel systeem en een machine kan de stellingen van een formeel systeem volgens de regels van dat systeem produceren. Vervolgens bespreek ik een specifiek, in dit verband belangrijk, formeel systeem: de Rekenkunde (2.3). De Gödelstellingen worden in (2.4) geïntroduceerd, vergezeld van een schets van hun bewijs. (2.5) bevat een verrassing: ik geef een voorbeeld van een niet-formeel (maar wel precies beschreven) systeem. We zullen zien dat we ons heel goed een machine kunnen voorstellen die stellingen produceert volgens de regels van dit systeem. Dit niet-formele systeem zal een rol spelen bij onze bespreking van de stappen (ii) en (iii) van Lucas' Argument.

2.1 Machines

Bij 'machine' in Lucas' Argument moeten we ons duidelijk niet een stoommachine of weefgetouw voorstellen, maar een geprogrammeerde digitale computer. De computer heeft een eindig aantal mogelijke discrete interne machinetoestanden en een in elk stadium eindig, maar in principe willekeurig uitbreidbaar geheugen.⁽¹⁰⁾ Bovendien is hij uitgerust met een invoerpoort en een uitvoerpoort. Via de invoerpoort kunnen we de momentane machinetoestand en/of de momentane geheugeninhoud wijzigen. Via de uitvoerpoort komt bepaalde geselecteerde informatie aangaande de interne toestand en de geheugeninhoud beschikbaar. Het rekenproces verloopt stapsgewijs. Een rekenstadium wordt volledig gegeven door de interne machinetoestand en de inhoud van het geheugen. Het programma specificereert een aantal regels om (we nemen hier aan: deterministisch) van rekenstadium naar rekenstadium te gaan. Bij het doen van zo'n rekenstap verandert de interne toestand van de machine en/of wordt de geheugeninhoud gewijzigd.

Als je aan een werkende computer denkt zijn er eigenlijk twee dingen: de abstracte machine - het ontwerp- en haar fysische realisatie. De fysische realisatie zal zich natuurlijk niet altijd conform het ontwerp gedragen: we spreken dan van een hardware-fout. We zullen zien dat de abstracte machine *gerepresenteerd* kan worden als een formeel systeem, de fysische machine is dan een *instantiatie* van dat systeem.

2.2 Formele systemen

Een formeel systeem wordt gegeven door drie dingen: zijn taal, zijn axioma's en zijn afleidingsregels.

2.2.1 De taal

De taal van het systeem wordt precies vastgelegd door een uitputtende beschrijving van de simpelste elementen van die taal (letters, haakjes, enzovoorts) en door instructies hoe je uit elementen van de taal die je al geconstrueerd hebt nieuwe kunt maken.

Een enigszins gesimplificeerd voorbeeld: We beschrijven de taal van een formeel systeem dat we MINI zullen dopen. De taal van MINI is als volgt gegeven:

- a) P zit in de taal.
- b) Als \square en \square in de taal zitten dan zit ook $(\square \& \square)$ erin.
- c) Niets zit in de taal dan op grond van (a) en (b) toegelaten is.

In de taal zitten bijvoorbeeld: P, (P&P), (P&(P&P)), ((P&P)&P), ((P&P)&(P&P)) ; maar niet:))PP&PP((, en ook niet: (P&Q) .

2.2.2 De axioma's

De axioma's zijn een aantal (eventueel oneindig veel) zinnen uit de taal die als uitgangspunt dienen van formele bewijzen. We eisen dat de axioma's *effectief opsombaar* zijn. Dat betekent dat een digitale computer in principe een (eventueel oneindig doorgroeiende) lijst moet kunnen produceren waar elk axioma op den duur in voorkomt.

2.2.3 De regels

De regels van een formeel systeem zijn instructies hoe je beginnende met de axioma's *formele bewijzen* kunt construeren. (We eisen net als bij de axioma's dat de regels effectief opsombaar zijn.) Een formeel bewijs is een eindig rijtje zinnen: $\square_1; \square_2; \square_3; \dots; \square_n$. Elk van deze zinnen is een zin uit de taal van het systeem. Een zin kan op twee manieren in het rijtje voorkomen: de zin is een axioma of hij is door toepassing van een regel uit voorafgaande zinnen in het rijtje te concluderen. \square_n is de *conclusie* van het bewijs. Een zin is *bewijsbaar in het systeem* als er een bewijs is met die zin als conclusie.

MINI heeft als axioma: $((P \& P) \& P)$ en als regel: uit $(\square \& \square)$ mag je \square concluderen. Een voorbeeld van een bewijs is nu: $((P \& P) \& P); (P \& P); P$

MINI heeft precies drie bewijzen. Welke zijn dit? Welke stellingen zijn bewijsbaar in MINI?

De begrippen 'formeel systeem' en 'formeel bewijs', zoals hier geïntroduceerd, zijn te beschouwen als puur *syntactische* noties. Syntaxis gaat over allerlei eigenschappen van talen die niet te maken hebben met betekenis. Tot het domein van de Syntaxis behoren operaties op talige entiteiten, zoals letters, woorden, zinnen.

2.2.4 Het formele systeem van een machine

Zoals we gezien hebben, gaat een machine stapsgewijs van rekenstadium naar rekenstadium. Een rekenstadium wordt gegeven door: één uit een eindig aantal interne toestanden, plus de geheugen-inhoud. Hierboven werd reeds beschreven dat de geheugentoeestand in elk stadium eindig is. Een rekenstadium is dus per definitie eindig. We kunnen zulke rekenstadia representeren door zinnen van een formeel systeem. Op deze manier kunnen we de *abstracte* machine representeren als formeel systeem. Als de machine bijvoorbeeld drie interne

toestanden: S_1, S_2, S_3 heeft en de mogelijke geheugeninhouden bestaan uit rijen nullen en enen, dan zouden: $S_1;0010$, $S_2;01110$, en $S_3;111$ zulke zinnen kunnen zijn.

We kunnen de rekenregels van de machine vertalen naar afleidingsregels van ons systeem. Het beginstadium van het rekenproces laten we corresponderen met één uniek axioma. Bewijzen in ons systeem zullen overeenkomen met rijtjes van opeenvolgende rekenstadia vanaf de begintoestand. Mocht de machine stoppen, dan wordt de hele rekenprocedure beschreven door het langste bewijs in het systeem.

We hebben nu de machine abstract gerepresenteerd als formeel systeem, zodat we de *fysische* machine kunnen zien als instantiatie of belichaming van dit systeem.

Merk op de formele systemen die bij machines horen heel specifiek zijn: zo is er maar één axioma en kunnen bewijzen op hoogstens één manier met één zin verlengd worden.

2.2.5 Machines kunnen formele systemen opsommen

Behalve dat machines beschreven kunnen worden met behulp van een formeel systeem, kunnen ze ook worden gebruikt om zelf een formeel systeem te produceren. De stellingen van een formeel systeem kunnen worden *opgesomd* door een geschikt geprogrammeerde machine. Opsommen betekent hier dat een reeks stellingen uit de uitvoerpoort van de machine zal komen, waarbij elke stelling van het systeem ooit zal verschijnen. Ik zal zeggen dat de machine het formele systeem opsomt als zij niet alleen de stellingen van het systeem produceert, maar dit ook doet *volgens de regels van het systeem*; ik bedoel hiermee dat de machine intern de axioma's van het systeem genereert, bewijzen bouwt uit die axioma's volgens de regels van het systeem en tenslotte de conclusies van deze bewijzen als uitvoer presenteert. Het hoeft echter niet zo te zijn dat de regels van het opgesomde systeem identiek zijn aan de regels volgens welke de machine werkzaam is (i.e. de regels van het formele systeem dat de machine abstract beschrijft zoals in 2.2.4 besproken). Het kan heel goed zijn dat de machine vele rekenstappen nodig heeft om één bewijsstap van het opgesomde systeem uit te voeren. Bovendien zal de machine de bewijzen van het formele systeem in een bepaalde volgorde moeten aflopen. De machine zal deze volgorde moeten bepalen. De instructies voor het bepalen van deze volgorde zullen in het algemeen geen deel uitmaken van het formele systeem dat door de machine wordt geproduceerd.

2.2.6 Interpretatie

De systemen waarvoor de Gödelstellingen geformuleerd zijn, zijn *interpreteerbaar*. We kunnen aan de symbolen van deze systemen betekenissen toekennen, zodat de zinnen van het systeem -gegeven deze betekenissen- opgevat kunnen worden als *uitspraken*. De verleende betekenissen maken de taal van het systeem pas echt tot *taal* in de gebruikelijke zin des woords: namelijk iets waarmee je iets kunt zeggen.

Bezie nu een systeem, waarvoor een interpretatie gespecificeerd is. We willen natuurlijk graag dat de axioma's *waar* zijn. (Een zin is waar als wat die zin zegt zo is; "Er staat een raket naast de Dom", bijvoorbeeld, is precies dan waar als er daadwerkelijk een raket naast de Dom staat. 'Waarheid' is ook van toepassing op de taal van ons systeem, nu zijn zinnen onder de gegeven interpretatie iets uitdrukken.) Verder is het gewenst dat de regels van het systeem *correct* zijn. Een regel is correct wanneer hij bij ware premissen altijd een ware conclusie oplevert. Omdat het soms erg moeilijk is om in te zien of zekere voorgestelde axioma's waar zijn, behouden we ons het recht voor af en toe geïnterpreteerde systemen met eventueel onware axioma's te beschouwen.

We interpreteren de taal van MINI: laat "P" staan voor: De Dom staat in Utrecht. Laat "&" staan voor: en. Nu wordt "(P&P)" geïnterpreteerd als: De Dom staat in Utrecht en de Dom staat in Utrecht. Merk op dat de haakjes geen interpretatie krijgen: ze functioneren slechts als hulpsymbolen.

Ik geef toe: geen erg inspirerend voorbeeld. De lezer overtuige zichzelf van de waarheid van het axioma en de correctheid van de enige regel van MINI.

2.2.7 Rijkdom

De formele systemen waarvoor de Gödelstellingen geformuleerd zijn, bevatten doorgaans de axioma's en regels van de zogeheten *eerste orde predikatenlogica*, plus de axioma's van de Elementaire Rekenkunde. We zullen zulke systemen *rijk* noemen. De eerste orde predikatenlogica codificeert het gewone wiskundige redeneren: zij bevat de regels voor het gebruik van uitdrukkingen als: en, of, niet, als ... dan, voor alle, er is een. Systemen in de eerste orde predikatenlogica hebben oneindig veel bewijzen; de totaliteit van deze bewijzen is in het algemeen zeer gecompliceerd. De Rekenkunde wordt behandeld in 2.3

2.2.8 Consistentie

Een formeel systeem heet *consistent* als niet elke zin uit de taal van het systeem bewezen kan worden. Als een systeem de eerste orde predikaten logica omvat kunnen we ook een andere definitie geven: het systeem is consistent als er geen zin \square is, zodanig dat zowel \square als de ontkenning van \square bewezen kunnen worden. Of nog anders: als we geen *tegenspraak* van de vorm ' \square en niet- \square ', kunnen bewijzen. Hoe kunnen we erachter komen of een gegeven systeem consistent is? Helaas, er is geen algemene procedure om de consistentie te bepalen. Bezie maar eens een formeel systeem met oneindig veel bewijzen; consistentie voor zo'n systeem betekent dat er een zin \square is zodat er geen bewijs is van \square . We moeten dus van elk bewijs inzien dat het niet- \square als conclusie heeft. De naieve manier om hier achter te komen zou zijn de bewijzen één voor één te bekijken en te controleren dat hun conclusie niet- \square is. Omdat - naar we hebben aangenomen- er oneindig veel bewijzen zijn komt deze procedure, als \square inderdaad niet als conclusie voorkomt, niet tot een eind. Het begint er nu naar uit te zien dat we er nooit achter kunnen komen of een formeel systeem met oneindig veel bewijzen consistent is. Gelukkig kunnen we in veel gevallen andere procedures bedenken dan de naieve. Stel, bijvoorbeeld, je hebt een geïnterpreteerd systeem in de eerste orde predikatenlogica. Stel ook dat we kunnen inzien dat de axioma's waar zijn en de regels correct. Dan is het niet moeilijk aan te tonen dat alle stellingen van het systeem waar zijn. Er volgt nu dat het systeem consistent is: als zowel \square als de negatie van \square bewijsbaar waren, zouden immers zowel \square als de negatie van \square waar zijn. Dit laatste is onmogelijk.

2.2.9 Formele bewijzen en echte bewijzen⁽¹¹⁾

Er is iets vreemds aan de hand met de notie 'bewijs' sinds Euclides. De begrippen 'axiomatisch bewijs' en, later, 'formeel bewijs' waren zo helder en zo succesvol dat men het alledaagse begrip als vaag, twijfelachtig en zelfs non-existent ging beschouwen.⁽¹²⁾ "Echte bewijzen vind je alleen in Wiskunde en Logica." Ik deel dit scepticisme niet, in feite geloof ik dat echte bewijzen even gewoon zijn en even veelvuldig voorkomen als fietsen, tafels, bomen enzovoorts. Het lijkt me dan ook zinnig je af te vragen of formele bewijzen bewijzen zijn.

Het is zeker niet zo dat alle formele bewijzen ook echt bewijzen zijn: bijvoorbeeld een bewijs in het formele systeem van een machine is zeker geen bewijs in enige gebruikelijke zin. Het lijkt me echter dat als we een geïnterpreteerd systeem in gedachten hebben en als we van de axioma's van dat systeem inzien dat ze waar zijn en als we van de regels inzien dat ze correct zijn dat we dan zeker kunnen zeggen dat de bewijzen van dat systeem ook echte bewijzen zijn.

Laat ik proberen het contrast tussen de begrippen 'formeel bewijs' en 'echt bewijs' zo scherp mogelijk te maken. Het is in het gewone begrip 'bewijs' ingebouwd dat als je een bewijs van iets hebt dat dat dan ook waar is. Bijvoorbeeld als ik een bewijs heb dat sneeuw wit is, dan is sneeuw ook wit. Maar geeft het hebben van een bewijs op die manier dan absolute zekerheid? Nee, we kunnen ons namelijk vergissen in die zin dat we abusievelijk denken dat we een bewijs hebben, terwijl datgene wat we voor een bewijs houden er geen is. Ik kan bijvoorbeeld denken dat ik een bewijs heb dat $733.721=508493$ -namelijk een berekening volgens het gewone algoritme voor decimale vermenigvuldiging-, terwijl ik in werkelijkheid een rekenfout gemaakt heb. In dat geval heb ik geen berekening volgens het gewone algoritme uitgevoerd en heb ik ook geen bewijs. Bezie nu eens een geïnterpreteerd formeel systeem, zeg S , en stel dat we de waarheid van de axioma's van S en de correctheid van de regels van S inzien. (De Rekenkunde (zie 2.3) is een voorbeeld van zo'n systeem.) Formele bewijzen in S zullen nu ook echte bewijzen zijn. Laten we nu ook nog aannemen dat elk bewijs dat we überhaupt in de taal van S kunnen formuleren correspondeert met een formeel bewijs van S . (Deze hypothese is zeker niet waar voor de Rekenkunde.) Ik beweer dat zelfs in dit opgeklopte geval de begrippen 'formeel bewijs in S ' en 'bewijs in de taal van S ' niet samenvallen. Mijn argument is dit: zelfs al zien we terecht in dat de axioma's waar zijn en de regels correct, we kunnen toch de mogelijkheid openlaten dat we ons vergissen. Als we ons vergissen, dan zijn de formele bewijzen van S nog steeds formele bewijzen van S , ze zijn echter geen bewijzen meer. We zien dus dat de formele bewijzen van S bewijzen zijn, maar dat het nog steeds voorstelbaar is dat ze het niet zijn. Ergo de begrippen 'formeel bewijs in S ' en 'bewijs in de taal van S ' vallen niet samen.

Het hele eiereten is dat 'waarheid' en 'bewijsbaarheid' begrippen zijn *in termen waarvan we cognitief succes formuleren*, terwijl een formeel systeem je een proces geeft waarmee je eventueel bewijzen kunt produceren. Vergelijk het met een stenenfabriek: aan de ene kant heb je de criteria waaraan een goede steen moet voldoen, aan de andere kant heb je een productieproces dat zulke goede stenen oplevert. Ook al levert het proces de facto goede stenen, nog steeds zijn 'goede steen' en 'steen geproduceerd volgens het proces' niet synoniem: we kunnen ons voorstellen dat om een of andere reden die we totaal over het hoofd hebben gezien het proces plotseling slechte stenen zou gaan opleveren. 'Goede steen' correspondeert hier met 'bewijs', 'proces' met 'formeel systeem'.

2.3 De Rekenkunde

De Rekenkunde (meestal: *Peano Rekenkunde* genoemd) is een formeel systeem in de eerste orde predikatenlogica met een vaste interpretatie. De Rekenkunde gaat over de (natuurlijke) getallen: $0, 1, 2, \dots$. Je kunt in de taal zulke dingen uitdrukken als: $7+5=12$; $7 \cdot 5=35$; voor alle getallen m en n : $m+n=n+m$; voor alle getallen n is er een priemgetal p , met: $p > n$. (Een priemgetal is een natuurlijk getal dat groter is dan 1 en dat alleen deelbaar is door 1 en door zichzelf.) In de Rekenkunde kun je bovenstaande zinnen ook bewijzen. Het is niet moeilijk in te zien dat de axioma's van de Rekenkunde waar zijn en de regels correct. Er volgt dat je in de Rekenkunde alleen maar ware uitspraken over getallen kunt bewijzen en dat de Rekenkunde consistent is.

Bekijk nu eens de volgende zin: voor alle getallen n zijn er priemgetallen p en q zodat $p-q=2$ en $q > n$ (het zogeheten Priemtweelingen Vermoeden). We hebben niet het flauwste idee of

deze zin waar is; we hebben ook geen idee of deze zin -aannemende dat hij waar is- bewijsbaar is. De enige manier die we tot nu toe hebben kunnen bedenken om erachter te komen of die zin bewijsbaar is, is de naieve procedure van het één voor één aflopen van alle bewijzen, en die procedure hoeft niet tot een eind te komen.

Kan het überhaupt voorkomen dat een rekenkundige zin waar is, en niet bewijsbaar in de Rekenkunde? Het blijkt als je naar voorbeelden van rekenkundige zinnen kijkt waarvan je toevallig weet dat ze waar zijn -zoals: elk getal kan op precies één manier ontbonden worden in priemfactoren- dat deze ook in de Rekenkunde bewijsbaar zijn. De eerste Gödelstelling vertelt ons dat deze empirie bedrieglijk is: er is een ware rekenkundige zin die niet in de Rekenkunde bewijsbaar is.

2.4 Gödel's Onvolledigheidsstellingen

Er zijn twee onvolledigheidsstellingen. We concentreren ons voorlopig op de eerste. Gödel's Eerste Onvolledigheidsstelling vertelt ons, dat voor elk consistent formeel systeem dat voldoende Rekenkunde omvat, m.a.w. dat rijk is (zie 2.2.7), er een zin is -de Gödelzin van dat systeem- die niet bewijsbaar is in dat systeem. Bovendien is onder zekere verdere condities de negatie van de Gödelzin ook niet bewijsbaar.⁽¹³⁾

Het bewijs van de Eerste Onvolledigheidsstelling is geïnspireerd door een paradox over echte bewijsbaarheid. Bekijk de volgende zin eens:

Deze zin is niet bewijsbaar.

Veronderstel dat die zin bewijsbaar is. Dan is hij waar, immers alleen ware dingen zijn bewijsbaar. Maar die zin beweert van zichzelf dat hij niet bewijsbaar is. Dus als hij waar is, dan is hij onbewijsbaar. Er volgt -nog steeds aannemende dat die zin bewijsbaar is- dat de zin niet bewijsbaar is. Een tegenspraak. Het uitgangspunt dat de zin bewijsbaar is, leidt tot een tegenspraak en kan dus niet waar zijn. We moeten concluderen dat die zin niet bewijsbaar is. Maar dat laatste is precies wat de zin zegt en daarmee hebben we de zin bewezen!! We hebben nu een onbewijsbare zin bewezen: een *paradox* !

Wat er mis is met bovenstaande redenering is moeilijk te zeggen. Filosofen zitten daar elkaar nog steeds over in de haren. Over die vraag wil ik het hier ook niet hebben: het gaat immers om de Gödelstelling. Het idee van het bewijs van de Eerste Gödelstelling is als volgt: we gaan bovenstaande paradox zo goed mogelijk nabouwen, maar dan voor formele bewijsbaarheid in systemen die voldoende Rekenkunde bevatten. Omdat formele bewijsbaarheid een syntactisch, en dus glashelder begrip is, zullen we niet echt tot een paradox komen: de tot een paradox leidende redenering zal op subtiële manier geblokkeerd blijken te zijn. De beste benadering van de tot een paradox leidende redenering geeft ons echter een fraaie bonus: de Eerste Onvolledigheidsstelling. Om de gedachten te bepalen nemen we als formeel systeem de Rekenkunde zelf.

Laten we het eerst eens naief proberen en leren van wat er niet gaat. Onze eerste poging om de paradoxale zin te formuleren wordt:

Deze zin is niet bewijsbaar in de Rekenkunde.

Het zij duidelijk dat deze zin niets interessants kan opleveren: hij is op flauwe manier waar. Het is namelijk helemaal geen zin uit de taal van de Rekenkunde en dus zeker niet bewijsbaar in de Rekenkunde. Er zijn in feite twee niet-rekenkundige ingrediënten: 'formele

bewijsbaarheid' is geen rekenkundig begrip maar een begrip uit de elementaire Syntaxis en het indexicale 'deze' komt niet voor in de rekenkundige taal.

De grote vondst van Gödel is het idee om bovenstaande zin te 'vertalen' in de taal van de Rekenkunde door de elementaire Syntaxis in die taal te vertalen.

Laat ik eerst aanduiden hoe we het gebruik van 'deze' kunnen vervangen door een toepassing van een elementair syntactische operatie.

2.4.1 De eliminatie van 'deze'

De truc die we hier presenteren is afkomstig van W.V.O. Quine. Bezie bijvoorbeeld:

Deze zin is leuk.

Quine's parafrase is als volgt:

*"geeft iets leuks wanneer je het achter zijn eigen aanhaling zet." geeft iets leuks
wanneer je het achter zijn eigen aanhaling zet.*

De aanhaling van een aantal woorden is wat je krijgt als je die woorden tussen aanhalingstekens zet. Wat moeten we volgens onze zin tussen aanhalingstekens zetten? De woorden die in die zin tussen aanhalingstekens gepresenteerd worden! Het resultaat van deze operatie is de zin zelf. De zin zegt dus van zichzelf dat hij leuk is. Precies het resultaat dat we wilden hebben. De beschreven operatie: een aantal symbolen achter hun eigen aanhaling zetten, is van elementair syntactische aard.

2.4.2 Hoe elementaire Syntaxis te vertalen naar Rekenkunde

Om de vertaling van elementaire Syntaxis naar Rekenkunde uit te leggen, beperk ik me tot een zeer klein deeltheoretje van de elementaire Syntaxis. Het is de theorie van reeksen van de letters A en B met de operatie concatenatie (= achter elkaar zetten). Voorbeelden van zulke reeksen zijn ABBA en AABAAB. De concatenatie van deze reeksen is: ABBAABAAB. We schrijven: $ABBA * AABAAB = ABBAABAAB$. In de theorie hebben ook nog λ , de lege tekenreeks, en variabelen x, y, z voor willekeurige tekenreeksen.

Het symbool ' λ ' is een *naam* voor de lege tekenreeks; het is natuurlijk niet *zelf* de lege tekenreeks. Dit is wat verwarrend omdat bijvoorbeeld 'ABBA' niet slechts fungeert als naam van de tekenreeks ABBA, maar ook als voorbeeld van die tekenreeks zelf.

Enige kenmerken van concatenatie zijn:

- _ $(x * \lambda) = (\lambda * x) = x$.
 - _ $((x * y) * z) = (x * (y * z))$.
 - _ Als $x * y = x * z$, dan $y = z$.
-

We gaan nu ons theorieetje vertalen naar de Rekenkunde. Dit doen we door aan A, B, λ getallen toe te kennen en concatenatie te representeren door een rekenkundige operatie:

$$\begin{array}{l} \square \square 0 \\ A \square 1 \\ B \square 2 \end{array}$$

De operatie \square is als volgt gedefinieerd:

$m \square n := m \cdot (\text{de kleinste macht van } 2 \text{ groter dan } n) + n$. Bijvoorbeeld:

$$7 \square 5 = 7 \cdot (\text{de kleinste macht van } 2 \text{ groter dan } 5) + 5 = 7 \cdot 8 + 5 = 61.$$

We representeren nu \square door \square .

Wil \square de concatenatie geschikt representeren dan moet \square de verschillende eigenschappen van concatenatie hebben en moet dit feit in de Rekenkunde verifiëerbaar zijn. Bijvoorbeeld: $(7 \square 0) = (7 \cdot 1 + 0) = 7$, en $((3 \square 4) \square 5) = ((3 \cdot 8 + 4) \cdot 8 + 5) = 229 = (3 \cdot 64 + (4 \cdot 8 + 5)) = (3 \square (4 \square 5))$.

We vertalen nu ABBA, dat is $A \square B \square B \square A$, als: $1 \square 2 \square 2 \square 1$, dat is: 53.

Vertaling zoals we dat hier ingevoerd hebben heet in wiskundige taal: inbedding in een definitionele uitbreiding. Een belangrijk kenmerk hiervan wordt gegeven door het volgende plaatje:

$$\begin{array}{l} AB \square BA = ABBA \\ \square \square \square = \square \\ 6 \square 5 = 53 \end{array}$$

Einde van 2.4.2

We hebben nu de ingrediënten bijeen voor het bewijs van de Eerste Onvolledigheidsstelling: we hebben laten zien hoe we het indexicale 'deze' kunnen elimineren en hoe we de elementaire Syntaxis kunnen vertalen naar de Rekenkunde. We kunnen nu een rekenkundige zin construeren die hetzelfde effect heeft als: "*Deze zin is niet bewijsbaar in de Rekenkunde*". Laten we deze zin "G" dopen. Deze zin is natuurlijk niet paradoxaal: het zou toch te gek zijn als een puur rekenkundige zin paradoxaal was.⁽¹⁴⁾

We redeneren nu zo: stel G is formeel bewijsbaar in de Rekenkunde, dan is G ook waar, immers de axioma's van de Rekenkunde zijn waar en de regels zijn correct. Maar als G waar is, dan is G niet formeel bewijsbaar in de Rekenkunde. Dus, gesteld dat G formeel bewijsbaar is in de Rekenkunde, volgt dat G niet formeel bewijsbaar is in de Rekenkunde. Een tegenspraak. Er volgt dat G niet formeel bewijsbaar kan zijn in de Rekenkunde. Dit is echter weer wat G zegt: dus G is waar.

Kunnen we niet verder gaan en een paradox afleiden door bovenstaande redenering om te zetten in een redenering in de Rekenkunde en daarmee G formeel te bewijzen in de Rekenkunde? Nee dat gaat niet: de stap 'Als G formeel bewijsbaar is in de Rekenkunde, dan is G waar.' is niet verifiëerbaar in de Rekenkunde.

We hebben gezien dat G waar is en niet bewijsbaar in de Rekenkunde. De negatie van G is onwaar en daarom ook niet bewijsbaar in de Rekenkunde. Hiermee hebben we de Eerste Onvolledigheidsstelling van Gödel bewezen voor het geval van de Rekenkunde. In feite hebben we bij ons bewijs gebruik gemaakt van tamelijk zwaar geschut: we hebben het principe gebruikt dat wat bewijsbaar is in de Rekenkunde waar is. Het is mogelijk het bewijs te verfijnen zodat de enige substantiële onderstelling die je nodig hebt de consistentie van de Rekenkunde is.

Je kunt laten zien dat de stelling ook geldt voor iedere consistente, rijke theorie (zie 2.2.7). De eis van 'rijkdom' (of een soortgelijke eis) is nodig om de vertaling van Syntaxis in de beschouwde theorie mogelijk te maken.

We hebben gezien dat machines formele systemen kunnen opsommen. In feite kun je de Gödelzinnen van opgesomde systemen op volledig mechanische manier construeren uit programma's van de opsommende machines, mits zulke programma's gegeven zijn in een vast gekozen programmeertaal.

In 2.2.8 hebben we het begrip 'consistentie' ingevoerd. 'Consistentie' is een syntactische notie. We kunnen het consistentiebegrrip dus vertalen in de taal van de Rekenkunde. De Tweede Onvolledigheidsstelling vertelt ons dat als een consistente theorie T rijk is, dat dan T de vertaling in rekenkundige taal van " T is consistent ." niet bewijst. In slogan (maar minder precies): een voldoende rijke theorie kan haar eigen consistentie niet bewijzen, tenzij die theorie inconsistent is. Het bewijs van de Tweede Onvolledigheidsstelling is een verfijning van het bewijs van de Eerste. Het voert te ver om het hier te reproduceren.

Bekijk een theorie T die rijk genoeg is. Het gekke is dat als T bewijst dat zij consistent is, zij inconsistent *moet* zijn. Als T daarentegen bewijst dat zij inconsistent is, dan is het mogelijk dat T inderdaad inconsistent is, maar het zou ook kunnen dat T consistent is en (dus) onwaar.

2.5 De Feferman Rekenkunde⁽¹⁵⁾

Bewijzen in een formeel systeem vormen om vele redenen geen goed model van het menselijk redeneren. Eén punt is dat wanneer een formeel systeem in de eerste orde predikatenlogica op een tegenspraak stuit, het hele systeem 'explodeert' en alles bewijsbaar wordt. Wij mensen concluderen wanneer we op een tegenspraak stuiten natuurlijk niet dat alles waar is, maar we onderzoeken de relevante hypothesen en verwerpen een paar van de meest verdachte.

Bezie de gewone (Peano) Rekenkunde. Een eerste benadering van de boven beschreven procedure van het terugtrekken van hypothesen voor het geval van de Rekenkunde, verkrijgen we door de Peano Rekenkunde om te zetten in de Feferman Rekenkunde. Dat gaat zo. We beginnen eerst alle bewijzen van Peano Rekenkunde op een rij te zetten: $\square_1, \square_2, \square_3, \square_4, \dots$. Zodra we een bewijs van een tegenspraak tegenkomen, kijken we wat het *langste* rekenkundige axioma \square is in dat bewijs. We gooien nu alle bewijzen weg die gebruik maken van \square of van langere rekenkundige axioma's. Daarna gaan we door met opsommen, waarbij we alleen bewijzen accepteren die gebruik maken van rekenkundige axioma's korter dan \square . Als we eventueel weer een bewijs van een tegenspraak tegenkomen, zoeken we weer het langste rekenkundige axioma van dat bewijs en herhalen de procedure. Een zin heet nu

bewijsbaar in Feferman Rekenkunde als er bij de hierboven beschreven procedure *ooit* een bewijs van die zin wordt geproduceerd en dat bewijs *nooit* wordt weggegooid.

Feferman Rekenkunde is niet een formeel systeem in de gebruikelijke zin. Toch zou het misleidend zijn het een 'informeel' systeem te noemen: 'informeel' suggereert een gebrek aan precisie bij de beschrijving van het systeem; de bovenstaande beschrijving laat echter niets te wensen over.

De procedure van de Feferman Rekenkunde lijkt mij volstrekt redelijk. Het modelleert het gedrag van iemand die de axioma's en regels van de Peano Rekenkunde plausibel vindt, maar toch nog een slag om de arm houdt. Natuurlijk is de ingebouwde 'back tracking' procedure niet die welke een mens echt zou volgen: het verwerpen van axioma's gebeurt niet op grond van inzicht in hun onjuistheid, maar op grond van een puur formeel criterium. Toch is de Feferman Rekenkunde al een opmerkelijk nette theorie. Het is bijvoorbeeld eenvoudig in te zien dat zij een theorie is in de eerste orde predikatenlogica.

Wij, die weten dat de Rekenkunde consistent is, kunnen inzien dat een rekenkundige zin bewijsbaar is in Fefermans systeem precies dan als hij bewijsbaar is in de gewone Rekenkunde. We zullen immers nooit een bewijs van een tegenspraak tegenkomen bij het opsommen van de bewijzen van de Rekenkunde, dus zal de in Fefermans systeem ingebouwde 'back tracking' procedure nooit worden geactiveerd. Toch zijn Fefermans systeem en dat van Peano naar intentie verschillend.

We kunnen het begrip 'bewijsbaarheid in Feferman Rekenkunde' vertalen in de taal van de Rekenkunde. Dit stelt ons in staat het bewijs van Gödel's Eerste Onvolledigheidsstelling te herhalen voor de Feferman Rekenkunde.⁽¹⁶⁾ Wat we nodig hebben om de waarheid, en daarmee de niet bewijsbaarheid, van de Gödelzin van Fefermans systeem te bewijzen is niet de consistentie van Feferman Rekenkunde -dat is een te zwakke bewering: de consistentie is immers 'ingebouwd'- , maar de consistentie van Peano Rekenkunde.

De Tweede Onvolledigheidsstelling geldt niet voor Fefermans systeem. Wij kunnen zelf eenvoudig inzien dat Fefermans systeem consistent is, zelfs als we ons agnostisch opstellen ten aanzien van de consistentie van de Peano Rekenkunde. De eliminatie van tegenspraken is immers in het systeem ingebouwd! Het blijkt dat de redenering die aan dit inzicht ten grondslag ligt in de Rekenkunde (en dus ook in de Feferman Rekenkunde) vertaald kan worden. Dus bewijst de Feferman Rekenkunde haar eigen consistentie.

Het is duidelijk dat Fefermans systeem door een machine 'opgesomd' kan worden. De machine volgt gewoon de hierboven beschreven procedure. Laten we ons de uitvoerpoort van de machine voorstellen als een scherm, waarop de voorlopig bewezen stellingen één voor één getoond worden. Mocht de machine op een tegenspraak stuiten, dan verschijnt "Ik herroep de volgende stellingen:" gevolgd door een lijst van de conclusies van de weggegooide bewijzen op het scherm. Daarna vertelt de machine ons: "Ik ga door met voorlopige stellingen" en komen er weer voorlopige stellingen op het scherm. Als we onwetendheid pretenderen aangaande de consistentie van de gewone Rekenkunde, is het echter niet duidelijk welke stellingen die op het scherm verschijnen blijvend zijn en dus tot de uiteindelijke stellingen van het systeem behoren.⁽¹⁷⁾ Naast de machine die Fefermans systeem opsomt zou een machine kunnen staan die de gewone (Peano) Rekenkunde opsomt. We kunnen ons voorstellen dat de twee machines er van buiten precies hetzelfde uitzien en dat gelijktijdig dezelfde stellingen op beide schermen verschijnen. Omdat de Rekenkunde in feite consistent is zijn beide machines dus vanuit behaviouristisch oogpunt hetzelfde. In werkelijkheid volgen ze echter verschillende procedures!

§3 **Machina abscondita**⁽¹⁸⁾

We zullen de schematische versie van Lucas' Argument nog eens ten tonele voeren, om deze vervolgens aan te vullen en te bespreken.

- (i) Elke machine is een instantiatie (=belichaming) van een formeel systeem .
- (ii) Dus gegeven een machine die consistent is en in staat is om eenvoudige rekenkunde te doen, is er -volgens de stellingen van Gödel- een ware rekenkundige zin die de machine niet zal produceren, namelijk de Gödelzin van het systeem dat hoort bij de machine.
- (iii) Wij kunnen echter de waarheid van die Gödelzin inzien.
- (iv) Ergo machines vormen geen adequaat model van de menselijke geest.

Hoe mooi zou het niet zijn als dit argument correct was. We zouden dan een eenvoudige weerlegging van het Mechanisme hebben die noch berust op een taalkundige goocheltruc (à la: een robot kan geen gevoelens hebben, want per definitie is 'gevoelens hebben' alleen van toepassing op organismen, en robots zijn geen organismen) of op de metafysica van het lekkere gevoel in je maag -hoewel het belang van dat lekkere gevoel niet onderschat mag worden. Helaas ...

Ad (i)

Machines zijn hier blijkbaar geprogrammeerde digitale computers. De bewering uit (i) heb ik al uitvoerig besproken in 2.1 en 2.2.1 t/m 2.2.4 .

Ad (ii)

De tweede stap begint met "Dus gegeven ... ". Dat "gegeven" is een belangrijk punt in de discussie. Ik wil de bespreking ervan echter liever uitstellen tot we stap (iii) behandelen.

Verder is er sprake van een machine die consistent is en in staat is elementaire rekenkunde te doen. Maar wat betekent dat in hemelsnaam? Misschien wel dit: de machine produceert een groeiende lijst van rekenkundige beweringen. Deze beweringen moeten consistent zijn (in een nader te beschrijven zin) en de stellingen van de Rekenkunde omvatten. We moeten elke bewering zien als een definitief geponeerde stelling die niet tijdens het verdere rekenverloop kan worden ingetrokken. Wat maakt dat de machine 'definitieve beweringen doet' en niet bijvoorbeeld zomaar wat rekenkundige zinnen opsomt, of misschien bezig is rekenkundige zinnen te *ontkennen* , of een 'back tracking' procedure zoals bij Fefermans systeem uitvoert? Het antwoord moet waarschijnlijk luiden: de bedoelingen van de programmeur of de specificaties van het programma. Of is het misschien voldoende dat wij -als het ware ad hoc- de activiteit van de machine interpreteren als het poneren van stellingen? In dat geval moeten we zeggen dat de machine een onbegrensde rij symboolreeksen produceert en dat wij het zijn die aan deze activiteit betekenis toekennen. Het nadeel van beide manieren van zien is dat ze de discussie in het nadeel van de machine lijken te prejudiciëren: misschien wil de machine zèlf wel iets rekenkundigs poneren. Ik ben geneigd te denken dat de bedoelingen-van-de-programmeur/programmaspecificaties theorie het meest in de goede richting zit. (Voor deze theorie is er natuurlijk het probleem dat een programma lang niet altijd doet wat de programmeur wil.) Voor onze doeleinden is het niet zo belangrijk van welke theorie we uitgaan. Aangezien de 'ad-hoc'- interpretatie theorie in het voordeel van Lucas lijkt te werken zal ik van deze uitgaan.

Laten we aannemen dat de machine (definitieve) rekenkundige stellingen opsomt. Is dit voldoende? De formulering "in staat is eenvoudige rekenkunde te doen" suggereert dat de manier waarop de machine tot haar beweringen komt lijkt op het produceren van rekenkundige bewijzen. Die veronderstelling zou -denk ik- de toelaatbare programma's wel erg beperken. Bovendien: is het wel zo duidelijk dat de manier waarop een menselijke wiskundige tot een rekenkundige stelling komt lijkt op het produceren van een rekenkundig bewijs?

Nu beweert Lucas: het formele systeem dat bij een programma hoort heeft een Gödelzin. Hier begaat Lucas echter een elementaire fout: de Eerste Onvolledigheidsstelling is -tenminste in zijn gebruikelijke formulering- van toepassing op 'rijke' formele systemen, dat wil zeggen formele systemen in de eerste orde predikatenlogica die voldoende rekenkunde omvatten. *Het formele systeem van een programma is helemaal niet rijk*. Om één punt te noemen: in de predikatenlogica kan een uitspraak meer dan één direct gevolg hebben, terwijl in het formele systeem van een machine elke zin ten hoogste één opvolger heeft. (Zie verder 2.2.4 en 2.2.7.) Misschien bedoelt Lucas eigenlijk dat het *door de machine opgesomde* systeem een Gödelzin heeft. Maar het is nog helemaal niet duidelijk dat we de activiteit van de machine willen zien als het opsommen van een systeem: het opsommen van rekenkundige beweringen is niet hetzelfde als het opsommen van de bewijzen van een systeem die deze beweringen als conclusies hebben. Gelukkig is er een eenvoudige uitweg voor Lucas uit al deze problemen: we kunnen gewoon de door de machine opgesomde beweringen zien als *axioma's* van een systeem in de eerste orde predikatenlogica. We kunnen indien de machine niet de axioma's van de Rekenkunde opsomt deze zelf toevoegen, zodat het systeem dat we krijgen automatisch rijk wordt. Laten we het zo verkregen systeem voor het gemak '*het geassocieerde systeem*' noemen. Dit systeem zal eventueel meer stellingen hebben dan de door de machine opgesomde beweringen. We geven nu onze definitie van het begrip 'consistente machine': een machine is consistent als het geassocieerde systeem consistent is. Het geassocieerde systeem heeft een Gödelzin; deze zal -als de machine consistent is- a fortiori niet door de machine worden opgesomd.⁽¹⁹⁾

Ad (iii)

Kunnen wij de waarheid van de in stap (ii) te voorschijn getoverde Gödelzin inzien? Zò gesteld is de vraag volstrekt inhoudsloos: er valt niets te zeggen als we niet meer weten over de manier waarop de machine aan ons gepresenteerd wordt. De grote vraag is dus: wat betekent het "gegeven" uit stap (ii)? Er zijn meerdere interpretaties mogelijk. Lucas is het duidelijkst over welke interpretatie hem voor ogen staat in J.R. Lucas (1968). Ik behandel eerst de meest voor de hand liggende interpretatie van dat "gegeven", die niet Lucas' eigen interpretatie is, maar wel die van sommige commentatoren op J.R. Lucas (1961) - bijvoorbeeld van P. Benacerraf (1967):

Ik loop over de Oude Gracht. Plotseling komt er een glimmende metalen doos op wieltjes en met flikkerende gekleurde lampjes aanrijden. Nadere beschouwing van de doos leert me dat er aan de voorkant een scherm zit waarop met geregelde tussenpozen rekenkundige stellingen verschijnen: $6+7=13$, voor elk natuurlijk getal n is er een priemgetal p met $p>n$, $1+2+\dots+n=1/2.n.(n+1)$, Kan ik nu de Gödelzin van het geassocieerde systeem bepalen? Nee, hiervoor heb ik het programma van de machine nodig. Dat programma is immers essentieel voor de beschrijving van het geassocieerde systeem. We moeten dus ons scenario uitbreiden. Een breed grijnzende meneer komt op me toestappen. "Aardig machientje, niet?" Hij laat me het vele vellen beslaande programma zien. Nu kan ik de Gödelzin van het geassocieerde systeem bepalen -in ieder geval *in principe*: de kans is groot dat ik voor het opschrijven van de Gödelzin meer symbolen nodig heb dan er elementaire deeltjes zijn in het Universum Nu is de vraag: kan ik de waarheid van die Gödelzin bepalen? De Gödelzin is

waar precies dan als het geassocieerde systeem consistent is. Ik moet dus uitmaken of het systeem consistent is. Dit is echter in het algemeen niet mogelijk, want: (a) stel je bijvoorbeeld voor dat ik kan inzien dat de machine precies de axioma's van de Rekenkunde opsomt, plus: voor alle natuurlijke getallen n zijn er priemgetallen p en q zodat $p-q=2$ en $q>n$; het Priemtweelingen Vermoeden. Omdat ik geen flauw idee heb of de Rekenkunde deze laatste bewering weerlegt of niet, heb ik ook geen flauw idee of de geassocieerde theorie van deze machine consistent is of niet. Bovendien is het volstrekt onduidelijk of ik, aannemende dat de geassocieerde theorie consistent is, daar ooit achter kan komen. (b) Het is heel goed mogelijk dat ik -anders dan onder (a)- geen enkel overzicht heb van de opgesomde beweringen; misschien somt de machine braaf de stellingen van de Rekenkunde op, maar zal zij plotseling -zeg: over een miljoen jaar- '0=1' als output geven.⁽²⁰⁾

Zoals we gezien hebben is (iii) onder het naieve scenario niet te verdedigen. Laten we nu het scenario uit J.R. Lucas (1968) bezien. Volgens Lucas moeten we het "gegeven" interpreteren in de context van de discussie tussen de Mechanist en de antimechanist. Wat Lucas wil weerleggen is de bewering dat hij een machine *is*. Daartoe stelt hij zich voor dat er iemand naar hem toekomt met een geweldig dik pak papier en beweert: "Dit is het programma van de Lucas-machine" Nu kan Lucas concluderen dat de machine die door het programma wordt gespecificeerd consistent is. Het programma *is* volgens de Mechanist immers het programma van Lucas en Lucas weet dat in ieder geval *hijzelf* consistent is. Maar weet Lucas eigenlijk wel van zichzelf dat hij consistent is? Lucas zegt (in A.R. Anderson (1964), p53): "The fact that we are all sometimes inconsistent cannot be gainsaid, but from this it does not follow that we are tantamount to inconsistent systems. Our inconsistencies are mistakes rather than set policies. They correspond to the occasional malfunctioning of a machine, not to its normal scheme of operations. Witness to this that we eschew inconsistencies when we recognize them for what they are. If we really were inconsistent machines, we should remain content with our inconsistencies, and would happily affirm both halves of a contradiction." Met andere woorden: mensen zijn in essentie zelfcorrigerende wezens en derhalve per definitie consistent.⁽²¹⁾ Lucas is dus consistent en zo moet ook -op grond van de door de Mechanist zelf verstrekte informatie- de door het programma van de Mechanist gedefiniëerde machine consistent zijn. Nu zijn we klaar: Lucas kan de bij het hem gepresenteerde programma horende Gödelzin bepalen. De machine die door het programma gespecificeerd wordt is consistent. Dus Lucas kan deze Gödelzin als waarasserteren, terwijl de machine deze niet zal kunnen opsommen. Ergo Lucas is niet die machine.

De these dat een zeker programma het programma van Lucas *is*, is in feite een absurd sterke mechanistische bewering. Bovendien richt de argumentatie in deze vorm zich tegen een strikte identiteitstheorie die waarschijnlijk op flauwe 'taalkundige' gronden al niet juist kan zijn. Lucas' argumentatie zou echter ook werken -even aannemende dat deze überhaupt werkt- als we ons zouden beperken tot de these dat een zekere machine *dezelfde rekenkundige capaciteiten heeft* als Lucas. Voor de weerlegging van deze zwakkere these kunnen we het Lucasiaanse scenario ietwat veranderen: nu spreekt de breed grijnzende meneer Lucas aan op de Oude Gracht. Hij stelt zichzelf voor als de Mechanist. Trots toont hij zijn machine en het bijbehorende programma. Hij beweert pertinent dat deze machine in rekenkundig opzicht hetzelfde presteert als Lucas. Lucas poneert op grond van de hierboven beschreven redenering de Gödelzin van de machine en weerlegt daarmee de boude bewering van de nu enigszins sip kijkende meneer.

In mijn bespreking wil ik het gewijzigde Lucasiaanse scenario aanhouden. Als je immers de zwakkere these weerlegt, weerleg je daarmee ook de sterkere. Bovendien denkt het prettiger als je je Lucas voorstelt geconfronteerd met een echte machine, dan wanneer je je Lucas' strikte identiteit met een machine moet indenken.

Het is duidelijk dat Lucas tot de waarheid van de Gödelzin van de machine concludeert *op grond van de door de Mechanist verstrekte informatie*. Het is echter zeer de vraag of Lucas wel het recht heeft de uitingen van de Mechanist als informatie te behandelen. Met andere woorden: heeft hij wel reden de Mechanist op zijn woord te geloven? Vreemd genoeg lijkt het argument dat we onder de loupe hebben self-defeating te zijn: Lucas concludeert dat de Gödelzin die hij uit het programma extraheert waar is op grond van de informatie van de Mechanist dat dat programma in rekenkundig opzicht hetzelfde presteert als Lucas. Maar uit het verdere argument blijkt dat de Mechanist ongelijk heeft. Lucas had dus van het begin af aan al geen reden om de Mechanist te vertrouwen, en daarmee ook geen reden om te concluderen dat de gevonden Gödelzin waar is. Wat betekent dit nu? Het betekent dat we de onderstelling dat Lucas de Mechanist op zijn woord mag geloven in het argument als hypothese moeten behandelen en dat deze hypothese tot een tegenspraak leidt. De conclusie van het argument is derhalve dat Lucas geen reden kan hebben om de Mechanist te vertrouwen. Wil dit zeggen dat de Mechanist dus ook onwaarheid spreekt? Niet per se: het is heel goed denkbaar dat we geen goede reden hebben iemand te vertrouwen, terwijl hij/zij toch de waarheid spreekt.

We komen nu tot wat ik zie als het cruciale bezwaar tegen Lucas' scenario. Lucas verslaat de machine op grond van de door de Mechanist verstrekte informatie. Maar dit is duidelijk geen *eerlijke* overwinning! De vergelijking is alleen fair als we zowel Lucas als de machine dezelfde informatie verschaffen. Maar als we de machine informatie geven kan het geassocieerde systeem veranderen⁽²²⁾ en hoeft de Gödelzin van het geassocieerde systeem zoals het eerst was niet hetzelfde te zijn als de Gödelzin van het geassocieerde systeem van de geïnformeerde machine. Het is heel goed denkbaar dat de veranderde machine net als Lucas de Gödelzin van het eerdere geassocieerde systeem zal kunnen produceren. Volgens dit bezwaar zijn (ii) en (iii) onder het Lucasiaanse scenario verward en verwarrend omdat daar ten onrechte machines worden behandeld alsof ze niet het vermogen zouden hebben informatie op te nemen, terwijl dit vermogen -terecht- wel aan mensen wordt toegeschreven. In het bijzonder is de Gödelzin waarvan sprake is in (ii) en (iii) de Gödelzin van de ongeïnformeerde machine, maar zijn de 'wij' uit (iii) wel geïnformeerd.

We kunnen de structuur van de argumentatie hier illustreren met een puur machinaal voorbeeld. Stel je twee machines voor: de één, zeg P, somt de Peano Rekenkunde op. De ander, zeg F, somt de Feferman Rekenkunde op. We zullen in dit voorbeeld niet naar de geassocieerde systemen kijken, maar rechtstreeks naar de opgesomde systemen. Zoals we gezien hebben produceren beide machines precies dezelfde stellingen. Beide machines hebben een invoerpoort via welke we rekenkundige beweringen kunnen inlezen. Deze beweringen worden dan aan de axioma's van de machines toegevoegd. F voert in dat geval de 'back tracking' procedure uit op de bewijzen die geconstrueerd worden met de nieuwe axioma's. De uitspraak "F somt de zelfde stellingen op als P." kan in de taal van de Rekenkunde vertaald worden. (Deze uitspraak is in feite equivalent met de bewering dat de Rekenkunde consistent is.) Geen van beide machines zal deze uitspraak kunnen bewijzen. Beide machines zullen precies hetzelfde reageren, wanneer we de uitspraak via de invoerpoort als axioma toevoegen. Natuurlijk zal F -de machine die in ons voorbeeld correspondeert met Lucas- de Gödelzin van de Rekenkunde kunnen bewijzen, als we aan F bovengenoemde bewering toevoegen. Maar het is evengoed zo dat P de Gödelzin van Fefermans systeem kan bewijzen, als we aan P de informatie geven dat F en P in hun onveranderde vorm dezelfde stellingen opsommen.

We hebben betoogd dat *als* Lucas in staat is op grond van de door de Mechanist verstrekte informatie de waarheid van de Gödelzin van de machine te bepalen, dat dan nog niet duidelijk is dat de machine *met die informatie* niet hetzelfde kan presteren. Dit is wat mij betreft voldoende om Lucas' Argument te ontkrachten. Het is echter toch interessant ons af te vragen of Lucas inderdaad de waarheid van de Gödelzin van de machine kan bepalen. Dit kan als we de machine opvatten als een apparaat dat stellingen definitief poneert. Maar waarom is dit nodig? Waarom kunnen we een machine niet zien als een zelfcorrigerend apparaat? Lucas probeert redenen te geven waarom dit niet zou kunnen (zie bijvoorbeeld A.R. Anderson (1964),53-55): mechanische zelfcorrectie procedures zijn volgens hem wel mogelijk, maar inherent onredelijk: als wij op een tegenspraak stuiten verwerpen we niet zomaar wat axioma's op grond van een vast selectiecriterium; nee, wij analyseren wat er aan de hand is en komen zonodig tot hele nieuwe axioma's, desnoods ook in een nieuwe taal. Ik denk dat Lucas hier gewoon de machinale mogelijkheden onderschat: het is volstrekt onduidelijk wat voor slimme correctieprocedures er niet mogelijk zijn. Je zou je zelfs kunnen voorstellen dat de machine haar zelfcorrectie-protocol bijstelt in het licht van nieuwverkrege informatie! ⁽²³⁾ (Merk op dat als Lucas hier gelijk zou hebben, we Lucas' Argument verder wel konden missen: mensen zijn op redelijke manier zelfcorrigerend, machines niet, dus mensen zijn geen machines.)

Laten we bezien wat er gebeurt met het Lucasiaanse scenario als we de door de Mechanist gepresenteerde machine opvatten als zelfcorrigerend. Laten we ons voorstellen dat de machine in principe alle gevolgen in de eerste orde predicaten logica van de door haar opgesomde beweringen controleert en dat zij als ze op een tegenspraak stuit bepaalde veronderstellingen -waaronder enige die essentieel zijn bij de afleiding van die tegenspraak- zal intrekken. We nemen nu als axioma's voor het geassocieerde systeem in nieuwe zin de door de machine opgesomde stellingen die nooit worden ingetrokken. Het is eenvoudig in te zien dat het geassocieerde systeem consistent zal zijn; we hebben de mededeling dat Lucas en de machine dezelfde stellingen opsommen helemaal niet nodig om tot die conclusie te komen. De waarheid van de Gödelzin van het geassocieerde systeem volgt echter niet uit de consistentie van de machine: immers dit systeem is geen formeel systeem in de gebruikelijke zin, omdat zijn axioma's de stabiele beweringen zijn van een zelfcorrigerende procedure. De Eerste Onvolledigheidsstelling van Gödel is in zijn gebruikelijke vorm alleen van toepassing op rijke formele systemen.⁽²⁴⁾ Lucas kan dus langs deze weg niet tot de waarheid van deze Gödelzin concluderen!

Bezie de machine F die de Feferman Rekenkunde opsomt. De vertaling in de taal van de Rekenkunde van: "De Feferman Rekenkunde is consistent.", zal door de machine opgesomd worden. De Gödelzin van Feferman Rekenkunde zal niet opgesomd worden. Om tot de waarheid van deze zin te concluderen is het inzicht in de consistentie van Fefermans systeem niet toereikend: we moeten weten dat Peano Rekenkunde consistent is!

De slotsom is: hoe het ook zij, of de door de Mechanist gepresenteerde machine nu als zelfcorrigerend gezien wordt of niet, het is niet duidelijk dat Lucas met de door de Mechanist gepresenteerde informatie meer kan doen dan de machine.

Ad (iv)

Aangezien (iii) in geen van de beschouwde scenario's is vol te houden, kunnen we niet tot (iv) concluderen.

Let wel: het enige wat we hebben aangetoond is dat een bepaald argument tegen het Mechanisme faalt. Er volgt natuurlijk niet daaruit dat het Mechanisme juist is.

Wat denk ik nu zelf? Is het Mechanisme in één of andere vorm waar? Ik kan me eerlijk gezegd niet aan het gevoel onttrekken dat we eigenlijk niet goed weten wat de verschillende mechanistische thesen betekenen. Het zou me niets verbazen als de verdere ontwikkeling van de Kunstmatige Intelligentie onze kijk op de inhoud van deze thesen ingrijpend verandert. Eén wel zeer problematisch aspect van het Mechanisme zoals dat in Lucas' Argument naar voren komt is het abstraherend afzien van beperkingen aan tijd en geheugenruimte. Ik meen dat wij *-tenminste in de context van de mens-machine problematiek-* niet goed begrijpen welke rol deze abstractie speelt.

Tot besluit, wil ik Kreisel het laatste woord laten: "[...] wide-eyed enthusiasts of A.I. or, better, of digital intelligence, and indignant critics never give thought to the possibility that we know enough already to refute their rhetoric ; but, presumably, not enough to settle any significant issue."⁽²⁶⁾

Noten

- 1) Gödel's eigen ideeën hieromtrent zijn te vinden in: H. Wang (1974), p. 324-328.
- 2) Ik geef graag toe dat deze bewering ietwat overdreven is. Men kan natuurlijk willekeurig veel "Lucas-argumenten" verzinnen. Het is onrealistisch van één argumentatiestrategie te verwachten dat zij alle denkbare varianten in één klap weerlegt. Filosofische discussies zijn nooit echt afgesloten.
- 3) Het is interessant de Gödelstellingen te vergelijken met Heisenberg's Onzekerheidsrelatie. Het is belangrijk om in te zien dat Heisenberg's resultaat betrekkelijk weinig impliceert aangaande, zeg, de menselijke vrijheid.
- 4) Het probleem hier is analoog met de vraag of de genetische code wel echt een code is. Wat is de gecodeerde boodschap? Wie of wat is de zender? Wie is de ontvanger? Het belangrijke punt is dat het er niet zoveel toe doet of de genetische code een code is of niet; waar het om gaat is dat het voor veel doeleinden verhelderend is de genetische code te zien als code.
- 5) 'Mechanisch' moet hier geïnterpreteerd worden als: gelijkenis vertonend met machinale processen. Niet als: machinaal. De claim dat het menselijk denken machinaal is, stuit op hetzelfde bezwaar als these 1.1.1.
- 6) Voor een leesbare introductie tot het Functionalisme, zie J.A.Fodor (1981).
- 7) Een problematisch aspect van deze doctrine is dat verschillende objecten op dezelfde tijd op dezelfde plaats kunnen zijn. Dit verschijnsel hoeft niet alleen op te treden bij objecten uit verschillende niveaus. Bijvoorbeeld in bepaalde stadia kunnen plant en stengel samenvallen, zonder dat daarmee in die stadia de plant de stengel *is*. (Dit voorbeeld is afkomstig van S.A.Kripke.) Verder moeten we oppassen voor orgiën van reïficatie. Neem nu de postbode. 'Postbode' is een *rol* van een mens. Het gaat niet aan een speciaal object te postuleren: de postbode, met onderliggende mens.
- 8) Ik hoop dat mijn gebruik van het woord 'gemedieerd' niet verwarrend is. De fysische processen zitten niet tussen de niet-fysische oorzaak en het niet-fysische gevolg in. Als ik bijvoorbeeld een boom zie, is er niet éerst de boom, dán allerlei tussenliggende fysische processen in ether, oog en brein en dán tenslotte de geheimzinnige gebeurtenis van het 'zien'. Op het niveau van het zien is er alleen maar dit: ik, de boom, mijn zien ervan. We kunnen spreken van directe waarneming. De *onderliggende* fysische processen maken deze waarneming mogelijk. Het idee van *tussenliggende* fysische

processen -we zouden dat 'de Mythe van het Medium' kunnen noemen- berust op een niveauverwarring.

- 9) Aan te bevelen literatuur: éénvoudige, vlot leesbare introducties zijn te vinden in: J.N. Crossley (1972), E. Nagel & J.R. Newman (1958), D.R. Hofstadter (1979). Een uitstekend leerboek is: G.S. Boolos & R. Jeffrey (1974). Het allermooiste -maar helaas een stuk moeilijker- is nog altijd het grandioze artikel S. Feferman (1960). Voor historisch geïnteresseerden is het aan te bevelen S. Feferman & alii (1986) te raadplegen.
- 10) In feite is wat we tot de interne toestanden rekenen en wat tot het geheugen een kwestie van conventie.
- 11) Deze sectie bevat een paar van mijn eigen ideeën over echte bewijzen. De meeste logici zullen mijn opinies over deze materie waarschijnlijk niet delen.
- 12) Zo zegt Lucas in J.R. Lucas (1968), p174: "Many people do speak colloquially of their being able to prove things a machine cannot prove, and it would be permissible to define an absolute provability to accomodate this locution. Permissible, but inexpedient. Provability has been construed by mathematical logicians for a generation as a syntactical term with a very precise definition, and it could be confusing to import loose notions of what I can see to be true into this well-disciplined and useful concept."
- 13) Uit de constructie zal blijken dat onder de gewone interpretatie de Gödelzin waar is en dus de negatie van de Gödelzin onwaar. Het dilemma is dus: iets waars niet bewijzen of iets onwaars bewijzen. Er bestaat een andere zin, de zogenaamde Rosserzin van het systeem, waarvan je kunt laten zien dat als het systeem dat we beschouwen rijk genoeg is en consistent, dat noch die zin, noch zijn negatie bewijsbaar zijn in het systeem. Er zijn nog veel wildere zinnen uitgevonden; deze vallen helaas buiten het bestek van deze voetnoot ...
- 14) De pointe is hier volgens mij dat 'echte bewijsbaarheid', net als 'waarheid', een *semantisch* begrip is. Van semantische begrippen zijn we paradoxen gewend. 'Formele bewijsbaarheid' is een *syntactisch* begrip: dat kan geen aanleiding geven tot paradoxen.
- 15) De Feferman Rekenkunde is uitgevonden door Feferman. (Dit is geen tautologie: de Peano Rekenkunde werd uitgevonden door Dedekind.) Zie S. Feferman (1960). De Feferman Rekenkunde is niet alleen filosofisch interessant: in de handen van Feferman en Orey werd zij een machtig hulpmiddel bij de wiskundige bestudering van Relatieve Interpreteerbaarheid. Mijn gebruik van de uitdrukking 'Feferman Rekenkunde' is trouwens idiosyncratisch: in de literatuur gaat het steeds over 'Feferman's Predicate'; hiermee bedoelt men: de vertaling in de taal van de Rekenkunde van het bewijsbaarheidspredicaat van Feferman Rekenkunde.
- 16) De Eerste Onvolledigheidsstelling voor Fefermans systeem volgt ook direct uit de Eerste Onvolledigheidsstelling voor de gewone Rekenkunde, gecombineerd met het inzicht dat de Feferman Rekenkunde en de Peano Rekenkunde de zelfde stellingen bewijzen. De Gödelzin van Feferman's systeem verschilt echter op interessante manier van die van Peano's systeem.
- 17) Wat kan de Feferman Rekenkunde over bewijsbaarheid in zichzelf zeggen? De situatie is nogal subtiel: voor *elke* rekenkundige zin \square geldt: als de Feferman Rekenkunde \square bewijst, dan bewijst de Feferman Rekenkunde ook dat de Feferman Rekenkunde \square bewijst. Dit betekent als het ware dat als de Feferman Rekenkunde iets bewijst, zij dit feit ook zelf kan inzien. Bezie echter de Gödelzin G van Fefermans systeem. Feferman Rekenkunde bewijst niet: als Feferman Rekenkunde G bewijst, dan bewijst Feferman Rekenkunde dat Feferman Rekenkunde G bewijst. Dit betekent dat de Feferman Rekenkunde niet zelf kan inzien dat als hij iets bewijst, dat hij dan ook bewijst dat hij dat iets bewijst.
- 18) Literatuur o.a.: P. Benacerraf (1967), G.L. Bowie (1982), D.R. Hofstadter (1979), J.R. Lucas (1961)&(1968), J.C. Webb (1983).

- 19) We hadden het probleem voor Lucas ook kunnen omzeilen door een andere metamathematische stelling te gebruiken dan de Gödelstelling. Voor de verdere discussie lijkt het er niets toe te doen welke oplossing we kiezen.
- 20) Een mogelijke tegenwerping tegen het argument zoals hier gepresenteerd is: het gaat er natuurlijk niet om wat ik hier en nu in kan zien, maar om wat ik in *principe* in kan zien, abstraherend van beperkingen van tijd, geheugenruimte enzovoorts. De principiële situatie voor het geval van de Rekenkunde plus het Priemtweelingen Vermoeden is als volgt: als de Rekenkunde plus het Priemtweelingen Vermoeden inconsistent is, dan kunnen wij daar in principe achter komen. Is deze theorie echter consistent, dan is het onduidelijk of wij daar zelfs in principe achter kunnen komen. De argumentatie onder (a), (b) kan overigens nog wat versterkt worden met metamathematische argumenten, zie bijvoorbeeld: G.L. Bowie (1982).
- 21) Ik denk overigens dat de vergelijking tussen "mistakes" en "the occasional malfunctioning of a machine" niet altijd opgaat: het lijkt nogal onwaarschijnlijk dat bijvoorbeeld Frege, die een inconsistent systeem bedacht, last had van een storing in zijn essentieel menselijke functioneren.
- 22) We zouden natuurlijk kunnen proberen te betogen, dat Lucas als hij nieuwe informatie krijgt dezelfde blijft, maar dat een werkende machine bij toevoer van informatie essentieel verandert, 'een andere machine wordt'. Dit lijkt mij weinig plausibel. Hoe het ook zij, als Lucas' Argument afhankelijk zou zijn van deze manoeuvre, dan zou het in feite niets anders zijn dan een flauwe taalkundige truc.
- 23) Een punt dat niet onvermeld mag blijven is dat menselijke zelfcorrectie gericht is op *waarheid*. Inconsistentie is alleen maar een indicator van onwaarheid. Een probleem in dit verband is dat de onwaarheid van een systeem -anders dan de inconsistentie- niet in principe herkenbaar hoeft te zijn.
- 24) Je kunt een voorbeeld geven van een zelfcorrigerend systeem (met een iets ingewikkelder procedure dan de in de text beschouwde⁽²⁵⁾) waarvan de stabiele beweringen gesloten zijn onder de regels van de predikaten logica en de stellingen van de Rekenkunde omvatten, terwijl de stabiele beweringen volledig zijn in de volgende zin: voor elke rekenkundige zin ϕ is hetzij ϕ hetzij de negatie van ϕ stabiel. Een gevolg is dat dit systeem meerdere Gödelzinnen zal hebben, waarvan sommige waar zijn, sommige onwaar. Dit betekent onder meer dat er een Gödelzin is die in het systeem zelf stabiel bewezen kan worden. Hieruit volgt niet dat het systeem inconsistent is -dat kan immers niet- maar alleen maar dat het een onware bewering stabiel bewijst. (Anders dan het Feferman systeem dat alleen maar voor de zekerheid een controle uitvoerde op bewijzen uit alleszins plausibele axioma's, is het hier genoemde systeem *experimenterend*: het probeert van elke zin uit of hij -in combinatie met zekere andere zinnen- consistent vol te houden is.)
- 25) Voor specialisten: de procedure komt neer op het zoeken van een oneindig pad in een recursieve binaire boom.
- 26) Deze opmerking is te vinden in een nog ongepubliceerd manuscript van G. Kreisel, getiteld: Gödel's excursions into intuitionistic logic.

Literatuur

- Anderson, A.R., *Minds and Machines*. Englewood Cliffs: Prentice-Hall, 1964.
 Benacerraf, P., God, the Devil, and Gödel. *The Monist* 51, 1967, 9-32.
 Boolos, G.S. & Jeffrey, R., *Computability and Logic*. Cambridge: Cambridge University Press, 1974.
 Bowie, G.L., Lucas number is finally up. *Journal of Philosophical logic* 11, 1982, 279-285.
 Crossley, J.N. & alii, *What is mathematical logic?*. Oxford: Oxford University Press, 1972.

- Feferman, S., Arithmetization of metamathematics in a general setting. *Fundamenta Mathematica XLIX*, 1960, 35-92.
- Feferman, S. & alii (eds.), *Kurt Gödel, Collected Works, Volume I, Publications 1929-1936*. Oxford: Oxford University Press, 1986.
- Fodor, J.A., The Mind-Body Problem. *Scientific American*, januari 1981, 124-132.
- Hofstadter, D.R., *Gödel, Escher, Bach*. New York: Basic Books, 1979. Nederlandse vertaling onder dezelfde titel: Amsterdam: Contact, 1985.
- Lucas, J.R., Minds, Machines and Gödel. *Philosophy* 36, 1961, 120-124. Herdrukt in: A.R. Anderson (1964), 43-59.
- Lucas, J.R., Satan stultified: a rejoinder to Paul Benacerraf. *The Monist* 52, 1968, 145-158.
- Nagel, E. & Newman, J.R., *Gödel's Proof*. New York: New York University Press, 1958. Nederlandse vertaling: *De stelling van Gödel*. Aula Paperback 136, 1986.
- Wang, H., *From Mathematics to Philosophy*. London: Routledge & Kegan Paul 1974.
- Webb, J.C., Gödel's Theorems and Church's Thesis. In: Cohen, R.S., and Wartofsky, M.W. (eds.), *Language, Logic, and Method*. Dordrecht: Reidel, 1983, 309-353.

Dit artikel is een bijdrage aan het septembernummer '*Geest, Computer, Kunst*' van het tijdschrift Grafiet. Ik dank de redacteuren van dit nummer, Peter Hagoort en Rob Maessen, voor hun gedetailleerde kritiek op eerdere versies. Ook mijn collega's Erik Krabbe en Fer-Jan de Vries stelden vele verbeteringen voor. Heb dank, Erik en Fer-Jan! Uiteraard ben alleen ik verantwoordelijk voor het eindprodukt.