

RCEC

Research Center voor

Examinering en Certificering

STANDAARDBEPALING TOETS GESPROKEN NEDERLANDS

**Research Center voor Examinering en Certificering
Ism Cito, Cinop consortium, TNO**

Inhoud

1. INLEIDING	5
1.1 ACHTERGROND EN ONDERZOEKSOPDRACHT.....	5
1.2 ONDERZOEKSAANPAK	6
1.3 LEESWIJZER.....	7
2. ONDERZOEK STANDAARDBEPALING TGN FASE 1	9
2.1 OPZET ONDERZOEK STANDAARDBEPALING TGN FASE 1.....	9
2.1.1 <i>Training en selectie van beoordelaars</i>	9
2.1.2 <i>Beoordeling van kandidaten</i>	10
2.1.3 <i>Selectie kandidaten onderzoek naar standaard A1min-niveau</i>	10
2.1.4 <i>Selectie kandidaten onderzoek naar standaard A1-niveau</i>	10
2.2 RESULTATEN ONDERZOEK STANDAARDBEPALING TGN FASE 1	11
2.2.1 <i>Resultaten Standaardbepaling A1min-niveau</i>	11
2.2.2 <i>Resultaten Standaardbepaling A1-niveau</i>	13
2.3 AANBEVELINGEN ONDERZOEK TGN FASE 2	15
3. ONDERZOEK STANDAARDBEPALING TGN FASE 2	17
3.1 OPZET ONDERZOEK STANDAARDBEPALING TGN FASE 2.....	17
3.1.1 <i>Training en selectie van beoordelaars</i>	17
3.1.2 <i>Beoordeling van kandidaten</i>	18
3.1.3 <i>Selectie kandidaten onderzoek naar standaard A1min-niveau</i>	18
3.1.4 <i>Selectie kandidaten onderzoek naar standaard A1-niveau</i>	18
3.1.5 <i>Selectie kandidaten onderzoek naar standaard A2-niveau</i>	19
3.2 RESULTATEN ONDERZOEK STANDAARDBEPALING TGN FASE 2	19
3.2.1 <i>Resultaten Standaardbepaling A1min-niveau</i>	19
3.2.2 <i>Resultaten Standaardbepaling A1-niveau</i>	20
3.2.3 <i>Resultaten Standaardbepaling A2-niveau</i>	21
4. CONCLUSIES EN AANBEVELINGEN	23
4.1 CONCLUSIES	23
4.2 AANBEVELINGEN.....	24
BIJLAGE 1 ONDERZOEK STANDAARDBEPALING TGN FASE 1	25
BIJLAGE 1.1 DATA TGN 2008 FASE 1, A1MIN-NIVEAU	27
BIJLAGE 1.2 DATA TGN 2008 FASE 1, A1-NIVEAU	29
BIJLAGE 2 ONDERZOEK STANDAARDBEPALING TGN FASE 2	31
BIJLAGE 2.1 DATA TGN 2008 FASE 2, A1MIN-NIVEAU.....	33
BIJLAGE 2.2.1 DATA TGN 2008 FASE 2, A1-NIVEAU DEEL 1	35
BIJLAGE 2.2.2 DATA TGN 2008 FASE 2, A1-NIVEAU DEEL 2	37
BIJLAGE 2.3 DATA TGN 2008 FASE 2, A2-NIVEAU	39
BIJLAGE 3 BETROUWBAARHEIDSINTERVAL STANDAARD A1MIN-NIVEAU.....	41

1. Inleiding

1.1 Achtergrond en onderzoeksopdracht

Het Ministerie van Justitie gaf in 2003 het Cinop consortium (Cinop, LTS en Ordinate) de opdracht een toets en afnamesysteem te ontwikkelen voor de toetsing van mondelinge vaardigheden in het Nederlands van buitenlanders. Volgens de opdracht zou de toets, die de Toets Gesproken Nederlands (TGN) genoemd zou gaan worden, in de eerste plaats geschikt moeten zijn voor de toetsing van de taalvaardigheid in het Nederlands van buitenlanders die zich duurzaam in Nederland wilden vestigen. Daarnaast zou de toets eventueel ook ingezet moeten kunnen worden in het kader van een examenstelsel voor inburgering in Nederland. Volgens de makers van de TGN meet de toets het gemak waarmee kandidaten in normaal conversatietempo gesproken Nederlands kunnen verstaan en begrijpen en hierop adequaat en verstaanbaar kunnen reageren.

De scoreschaal van de TGN heeft een ondergrens van 10 en een bovengrens van 80 scorepunten. De normering van de TGN is gerelateerd aan de niveaus van het Gemeenschappelijk Europees Referentiekader (CEF). Door de TGN worden de volgende cesuren en ranges van scorepunten voor de CEF-niveaus gehanteerd:

- 10-15 scorepunten is gedefinieerd als het CEF-niveau lager dan A1-min niveau;
- 16-25 scorepunten als het A1-min niveau;
- 26-36 scorepunten als het A1-niveau;
- 37-46 scorepunten als het A2-niveau;
- 47-56 scorepunten als het B1-niveau;
- 57-67 scorepunten als het B2-niveau;
- 68-79 scorepunten als het C1-niveau;
- 80 scorepunten als het C2-niveau.

Voor een gedetailleerde beschrijving van de TGN wordt verwezen naar het rapport 'Verantwoording Toets Gesproken Nederlands' (Cinop, 2005)

In 2006 en 2007 heeft TNO in overleg met het Ministerie van Justitie en het Cinop consortium onderzoek gedaan naar (onderzoeksvraag 1) of er substantiële verschillen in de beoordelingen waren tussen het systeem dat automatisch uitslagen genereert en menselijke beoordelingen, en (onderzoeksvraag 2) of de zak-slaaggrens, dwz de grens van 16 scorepunten om te kunnen slagen voor het A1-min niveau, van de TGN op het goede niveau ingesteld was. Met betrekking tot onderzoeksvraag 2 was de conclusie van TNO dat de zak-slaaggrens of cesuur voor het A1-min niveau te soepel ingesteld was. Voor een gedetailleerd verslag van het TNO onderzoek wordt verwezen naar het onderzoeksrapport 'Onderzoek naar de kwaliteit van het inburgeringsexamen buitenland.' (TN-DV-2007 C053).

In een debat dat plaatsvond op 5 juli 2007 gaf de Tweede Kamer aan dat zij naar aanleiding van het voornoemde onderzoek van TNO een extra onderzoek

wenste. Het Ministerie van VROM, Directoraat-Generaal Wonen, Wijken en Integratie, heeft het Research Center voor Certificering en Examinering (RCEC) daarop verzocht de volgende onderzoeksvraag te beantwoorden: 'Is het TNO-onderzoek naar de cesuurbepaling zorgvuldig uitgevoerd en is daarmee de conclusie gerechtvaardigd dat de cesuur voor de TGN te soepel is ingesteld?'

Op basis van studie van het TNO-onderzoek en eigen onderzoek bevestigde het RCEC de conclusie van het TNO-onderzoek dat de cesuur te soepel was ingesteld en concludeerde tevens dat op basis van de beschikbare data geen verantwoord antwoord gegeven kon worden op de vraag hoeveel strenger de cesuur zou moeten zijn. In de reactie van TNO op het RCEC-onderzoek werden bij die laatste conclusie kanttekeningen geplaatst. Voor een gedetailleerd verslag van het RCEC-onderzoek wordt verwezen naar het onderzoeksrapport 'Beoordeling TNO-rapport Onderzoek naar de kwaliteit van het inburgeringsexamen buitenland'. (RCEC, 2007).

In het onderzoeksrapport van het RCEC werden twee aanbevelingen gedaan. De eerste aanbeveling was om een nieuw onderzoek naar de cesuur op de TGN te laten uitvoeren. De tweede aanbeveling was om in afwachting van de uitkomsten van nieuw onderzoek de cesuur op de TGN aan te passen overeenkomstig het voorstel dat geformuleerd is in de brief van de minister voor Wonen, Wijken en Integratie van 29 mei aan de Tweede Kamer. Ook bij de tweede aanbeveling plaatste TNO in voornoemde reactie op het RCEC-onderzoek een aantal kanttekeningen. Als echter een pragmatische keuze voor een aangepaste cesuur gemaakt zou moeten worden, stelde TNO voor een cesuur tussen de oorspronkelijke cesuur van 16 scorepunten en de door het RCEC voorgestelde cesuur te hanteren. De minister voor Wonen, Wijken en Integratie heeft mede op basis van de uitkomsten van de diverse onderzoeken de oorspronkelijke cesuur met ingang van 15 maart 2008 tussentijds verhoogd. Deze tussentijdse verhoging ligt tussen de oorspronkelijke cesuur van 16 scorepunten en de cesuur die is voorgesteld per 29 mei (zie ook de brief van de minister voor Wonen, Wijken en Integratie van 28 november 2007).

Het Ministerie van VROM heeft de aanbeveling van het RCEC om nieuw onderzoek te doen naar het bepalen van de cesuur op de TGN overgenomen en begin 2008 het RCEC de opdracht gegeven om dit onderzoek te doen uitvoeren. In het onderhavige rapport wordt verslag gedaan van dit onderzoek.

1.2 Onderzoeksaanpak

Het doel van het onderzoek 'Standaardbepaling Toets Gesproken Nederlands' was te komen tot standaarden/cesuren voor onderscheiden niveaus op de Toets Gesproken Nederlands. Daartoe werd een onderzoek opgezet dat uit twee fasen bestond.

In Fase 1 wordt onderzoek naar de zogenaamde globale standaarden voor het

A1min-niveau en A1-niveau gedaan. Met behulp van de resultaten uit Fase 1 worden in Fase 2 de nauwkeurige standaarden voor het A1min-niveau en het A1-niveau bepaald. Bij aanvang van het onderzoek werd bepaald dat afhankelijk van de uitkomsten van Fase 1 het onderzoek van Fase 2 aangepast zou kunnen worden. Die aanpassing heeft inderdaad plaatsgevonden, dat wil zeggen dat het onderzoek uitgebreid werd met onderzoek naar de standaard voor het A2-niveau. De onderzoeksopzet van Fase 1 wordt in hoofdstuk 2 beschreven en de onderzoeksopzet van Fase 2 in hoofdstuk 3.

Het onderzoek is met de volgende organisaties uitgevoerd:

- RCEC: leiding en coördinatie van het project, analyseren van de gegevens en rapporteren.
- CINOP consortium: werving, training en selectie van beoordelaars, realiseren van de technische infrastructuur ten behoeve van de beoordeling van de kandidaten, verzameling van de beoordelingen, bijdragen aan onderzoeksopzet en rapportage.
- TNO: advisering bij de verschillende fasen van het project, bijdragen aan de onderzoeksopzet en de rapportage.

1.3 Leeswijzer

Hiervoor werden de achtergrond van het onderzoek, de onderzoeksopdracht en de onderzoeksaanpak beschreven. Hierna bestaat het rapport uit de volgende onderdelen. In hoofdstuk 2 wordt Fase 1 van het onderzoek naar de standaardbepaling van de TGN, dat wil zeggen het A1min-niveau en het A1-niveau, beschreven en worden de onderzoeksresultaten van Fase 1 gepresenteerd. In hoofdstuk 3 wordt Fase 2 van het onderzoek naar de standaardbepaling van de TGN, dat wil zeggen het A1min-niveau, het A1-niveau en het A2-niveau, beschreven en worden de onderzoeksresultaten van Fase 2 gepresenteerd. In hoofdstuk 4 worden conclusies gepresenteerd en aanbevelingen gedaan.

2. Onderzoek Standaardbepaling TGN Fase 1

Het onderzoek van Fase 1 heeft tot doel het bepalen van globale standaarden voor het A1min- en A1-niveau van de TGN. De uitkomsten van het onderzoek van Fase 1 zullen worden gebruikt voor het vaststellen van de onderzoeksopzet van Fase 2.

2.1 Opzet Onderzoek Standaardbepaling TGN Fase 1

2.1.1 Training en selectie van beoordelaars

De training die door kandidaat-beoordelaars gevolgd moest worden, was gericht op het verwerven of opfrissen van een goede kennis van het CEF. De training begon met een presentatie van het CEF. Achtereenvolgens kwamen de Algemene schalen, de schaal van Gesprekken Voeren en de Aspectschalen aan de orde. De CEF-niveaus en het Nederlandse A1-min niveau werden toegelicht met behulp van videomateriaal. Het videomateriaal laat gesprekken zien tussen een moedertaalspreker van het Nederlands en inburgeraars (NT2 sprekers). Er werd videomateriaal van gesprekken van kandidaten van verschillende niveaus gebruikt en tevens gewerkt met prompts uit het Basisexamen Inburgering. Omdat er geen authentiek materiaal, dwz examens van echte kandidaten, gebruikt mocht worden in de training, is tijdens de training gebruik gemaakt van voorbeelduitingen en videomateriaal. Kandidaat-beoordelaars kregen een prompt uit het examen te horen waarna de CINOP trainers uitgaande van CEF descriptorren de mogelijke responsen beschreven. Op deze manier werd de brug geslagen tussen het beoordelen op basis van de open gesprekken in het videomateriaal en de gesloten responsen in de uiteindelijke beoordelingstaak. .

Voor de beoordeling van kandidaten werden ervaren NT2 docenten geworven. Na het volgen van de training namen al de kandidaat-beoordelaars deel aan een assessment voor het A1min-niveau en een assessment voor het A1-niveau. Elk assessment bestond uit het beoordelen van 10 kandidaten. Van de 10 kandidaten die volgens hun TGN score tot het A1min-niveau behoorden, moesten de beoordelaars de volgende (dichotome) beslissing nemen:

- De kandidaat voldoet niet aan de beschrijving van het A1min-niveau of
 - De kandidaat voldoet aan de beschrijving van het A1min-niveau of hoger
- Van de 10 kandidaten die volgens hun TGN score tot het A1-niveau behoorden, moesten de beoordelaars de volgende (dichotome) beslissing nemen:
- De kandidaat voldoet niet aan de beschrijving van het A1-niveau of
 - De kandidaat voldoet aan de beschrijving van het A1-niveau of hoger.

Bij de selectie van kandidaat-beoordelaars werd gekeken naar de mate van onderlinge afwijking. Het selectie criterium was dat een kandidaat-beoordelaar die beoordelingen gaf die afweken van 80% van de kandidaat-beoordelaars afgewezen werd. Het oordeel van de trainers over de juistheid van de toegekende niveaus speelde geen rol bij de selectie van kandidaat-beoordelaars noch werden kandidaat-beoordelaars met een consistente afwijking, bijvoorbeeld milder of strenger dan de gemiddelde beoordelaar, bij voorbaat afgewezen. Kandidaat-beoordelaars kregen enkel te horen of zij wel of niet mochten deelnemen aan het onderzoek. De beoordelaars kregen geen inhoudelijke feedback op het gemaakte assessment om te voorkomen dat als het ware het CINOP consortium de standaarden voor de onderscheiden niveaus zou vaststellen. Geen van de kandidaat-beoordelaars bleek af te wijken van 80% van de kandidaat-beoordelaars. In totaal werden 25 beoordelaars geselecteerd.

2.1.2 Beoordeling van kandidaten

De beoordelaars kregen het examen van de kandidaat-inburgeraar te horen. De prompts waarop de kandidaat geacht werd te reageren, kregen de beoordelaars niet te horen maar wel te zien, zodat bij de beoordeling van de kandidaten ook de inhoudelijke correctheid van de responsen meegenomen kon worden. De beoordelingstaak voor de beoordelaars bestond uit het nemen van een (dichotome) beslissing, zie voorgaande paragraaf, of een kandidaat wel of niet voldeed aan de beschrijving van het A1min-niveau en het A1-niveau.

2.1.3 Selectie kandidaten onderzoek naar standaard A1min-niveau

Voor het onderzoek naar de globale standaard voor het A1min-niveau werden kandidaten geselecteerd die op de TGN de (door het consortium oorspronkelijk toegekende) scores 11, 14, 17, 20, 23, 26 en 29 behaald hadden. Deze scorepunten zijn evenredig verspreid over het deel van de TGN schaal waarbinnen naar verwachting de standaard voor het A1min-niveau ligt. Voor deze 7 scorepunten werden in totaal 21 kandidaten, 3 kandidaten per scorepunt, geselecteerd.

2.1.4 Selectie kandidaten onderzoek naar standaard A1-niveau

Voor het onderzoek naar de globale standaard voor het A1-niveau werden kandidaten geselecteerd die op de TGN de (door het consortium oorspronkelijk toegekende) scores 20, 23, 26, 29, 32, 35 en 38 behaald hadden. Deze scorepunten zijn evenredig verspreid over het deel van de TGN schaal waarbinnen naar verwachting de standaard voor het A1-niveau ligt. Voor deze 7 scorepunten werden in totaal 21 kandidaten, 3 kandidaten per scorepunt, geselecteerd.

2.2 Resultaten Onderzoek Standaardbepaling TGN Fase 1

2.2.1 Resultaten Standaardbepaling A1min-niveau

Bijlage 1.1 bevat de beoordelingen door 25 beoordelaars van de examens van 21 kandidaten. De beoordelaars moesten de volgende dichotome beslissing nemen:

- De kandidaat voldoet niet aan de beschrijving van het A1min-niveau of
- De kandidaat voldoet aan de beschrijving van het A1min-niveau of hoger

Met 3 kandidaten per scorepunt waarvan de examens door 25 beoordelaars beoordeeld worden, betekent het voorgaande dat er per scorepunt in totaal 75 examens dichotoom beoordeeld zijn. Van deze 75 examens is het percentage examens berekend dat als A1min-niveau of hoger beoordeeld wordt.

Tabel 2.1 bevat per scorepunt het percentage examens dat als A1min-niveau of hoger beoordeeld wordt.

Tabel 2.1: A1min-niveau
complete data

Score	Percentage
11	04
14	11
17	12
20	24
23	33
26	65
29	33

In Tabel 2.1 zou men met een toename van het aantal scorepunten een toename van het percentage examens mogen verwachten dat als A1min-niveau of hoger beoordeeld wordt. Als standaard voor het A1min-niveau zou dan in aanmerking kunnen komen het scorepunt waarvan ten minste 50% van de examens als A1min-niveau of hoger beoordeeld worden.

In Tabel 2.1 kan men lezen dat maximaal 33% van de examens van kandidaten met scorepunten 11 t/m 23 als A1min-niveau of hoger beoordeeld worden. De examens van de kandidaten met scorepunt 26 worden voor 65% als A1min-niveau of hoger aangemerkt. De verwachting dat een toename van het aantal scorepunten gepaard zou gaan met het percentage examens dat als A1min-niveau of hoger beoordeeld wordt, wordt echter door scorepunt 29 gelogenstraft. Op basis van de resultaten in Tabel 2.1 zou overwogen moeten worden om kandidaten met scorepunten 11 t/m 17 in Fase 2 van het onderzoek naar de standaard van het A1min-niveau buiten beschouwing te laten.

Om na te gaan of de resultaten van het onderzoek naar de standaard voor het A1min-niveau anders geweest zouden zijn indien zogenaamde afwijkende of

aberrante beoordelaars en/of kandidaten niet aan het onderzoek hadden deelgenomen, zijn een aantal heranalyses van de gegevens uitgevoerd. De resultaten van twee van die heranalyses worden hierna gepresenteerd.

De eerste heranalyse betrof het verwijderen van aberrante beoordelaars, dat wil zeggen de drie mildste beoordelaars (beoordelaars 18, 21 en 15) en de drie strengste beoordelaars (beoordelaars 14, 10 en 3). Tabel 2.2 bevat de resultaten van de eerste heranalyse.

Tabel 2.2: A1min-niveau
data exclusief 3 mildste en 3 strengste beoordelaars

Score	Percentage
11	02
14	07
17	07
20	19
23	33
26	68
29	32

Uit Tabel 2.2 blijkt dat de resultaten van de eerste heranalyse nauwelijks afwijken van de resultaten van de analyse van de complete data.

De tweede heranalyse betrof het verwijderen van verwijderde aberrante beoordelaars en aberrante kandidaten. Onder aberrante kandidaten worden hier kandidaten verstaan die beoordelingen van beoordelaars ontvangen die niet voldoen aan de verwachting dat een toename van het aantal scorepunten gepaard gaat met een toename van het percentage examens dat als A1min-niveau of hoger beoordeeld wordt. De kandidaten 4, 10, 15, 16 en 20 werden als aberrante kandidaten geïdentificeerd. Tabel 2.3 bevat de resultaten van de tweede heranalyse.

Tabel 2.3: A1min-niveau
data exclusief aberrante beoordelaars (strengste en mildste 3) en
exclusief aberrante kandidaten

Score	Percentage
11	02
14	00
17	07
20	29
23	03
26	89
29	37

Vergeleken met voorgaande analyses, blijkt het verwijderen van aberrante beoordelaars en kandidaten met name invloed te hebben op de percentages van de scorepunten 20, 23 en 26. Gegeven deze en voorgaande resultaten zou

overwogen kunnen worden om kandidaten met scorepunten 11 t/m 17 in Fase 2 van het onderzoek naar de standaard van het A1min-niveau buiten beschouwing te laten.

2.2.2 Resultaten Standaardbepaling A1-niveau

Bijlage 1.2 bevat de beoordelingen door 25 beoordelaars van de examens van 21 kandidaten. De beoordelaars moesten de volgende dichotome beslissing nemen:

- De kandidaat voldoet niet aan de beschrijving van het A1-niveau of
- De kandidaat voldoet aan de beschrijving van het A1-niveau of hoger

Met 3 kandidaten per scorepunt waarvan de examens door 25 beoordelaars beoordeeld worden, betekent het voorgaande dat er in totaal 75 examens dichotoom beoordeeld zijn. Van deze 75 examens is het percentage examens berekend dat als A1-niveau of hoger beoordeeld wordt.

Tabel 2.4 bevat per scorepunt het percentage examens dat als A1-niveau of hoger beoordeeld wordt.

Tabel 2.4: A1-niveau
complete data

Score	Percentage
20	01
23	25
26	05
29	23
32	27
35	31
38	67

In Tabel 2.4 valt het relatief hoge percentage bij scorepunt 23 en met name het hoge percentage bij scorepunt 38 op. Op basis van de resultaten in Tabel 2.4 zou overwogen kunnen worden om kandidaten met scorepunten 20 t/m 26 in Fase 2 van het onderzoek naar de standaard van het A1min-niveau buiten beschouwing te laten.

Ook voor het onderzoek naar de standaard voor het A1-niveau zijn een aantal heranalyses van de gegevens uitgevoerd waarvan de resultaten van twee heranalyses hierna worden gepresenteerd.

De eerste heranalyse betrof het verwijderen van aberrante beoordelaars, dat wil zeggen de drie mildste beoordelaars (beoordelaars 21, 13, 17) en de drie strengste beoordelaars (beoordelaars 3, 24, 14). Tabel 2.5 bevat de resultaten van de eerste heranalyse.

**Tabel 2.5: A1-niveau
data excl. aberrante (3 mildste en 3 strengste) beoordelaars**

Score	Percentage
20	02
23	23
26	05
29	21
32	23
35	32
38	72

Uit Tabel 2.5 blijkt dat de resultaten van de eerste heranalyse nauwelijks afwijken van de resultaten van de analyse van de complete data.

De tweede heranalyse betrof het verwijderen van aberrante beoordelaars en aberrante kandidaten (zie voorgaande paragraaf). De kandidaten 5, 12, 13 en 17 werden als aberrante kandidaten geïdentificeerd. Tabel 2.7 bevat de resultaten van de derde heranalyse.

**Tabel 2.6: A1-niveau
data excl. aberrante beoordelaars
(3 mildste en 3 strengste)
en excl. aberrante kandidaten**

Score	Percentage
20	02
23	13
26	05
29	03
32	32
35	11
38	72

Vergeleken met voorgaande analyses, blijkt het verwijderen van aberrante beoordelaars en kandidaten met name invloed te hebben op de percentages van de scorepunten 29 t/m 38. Gegeven de resultaten kan overwogen worden om kandidaten met scorepunten 20 t/m 26 in Fase 2 van het onderzoek naar de standaard van het A1-niveau buiten beschouwing te laten.

2.3 Aanbevelingen Onderzoek TGN Fase 2

Op basis van de resultaten van het onderzoek van Fase 1 naar de standaard voor het A1min-niveau en de standaard van het A1-niveau heeft de adviesgroep bestaande uit vertegenwoordigers van de bij het onderzoek betrokken organisaties (zie paragraaf 1.2) voor het onderzoek TGN Fase 2 besluiten genomen over:

- Het bereik van de te onderzoeken scorepunten voor het A1min-niveau en het A1-niveau;
- De (aanpassing van) de training en selectie van de beoordelaars;
- De (aanpassing van) taak van de beoordelaars;
- De selectie van het aantal beoordelaars;
- De selectie van het aantal kandidaten;
- De uitbreiding van het onderzoek met het A2-niveau.

Voor informatie over hoe voornoemde besluiten geconcretiseerd zijn, wordt verwezen naar hoofdstuk 3.

3. Onderzoek Standaardbepaling TGN Fase 2

In hoofdstuk 2 is verslag gedaan van Fase 1 van het onderzoek standaardbepaling TGN. In Fase 1 werd onderzoek gedaan naar de zogenaamde globale standaarden voor het A1min- en A1-niveau. Op basis van de resultaten van Fase 1, heeft de in hoofdstuk 2 genoemde adviesgroep besloten onderzoek te doen naar de standaarden voor de niveaus A1min, A1 en A2 volgens de onderzoeksopzet die in paragraaf 3.1 beschreven wordt. Na deze beschrijving worden de resultaten van de standaardbepaling voor de drie onderscheiden niveaus gepresenteerd en besproken.

3.1 Opzet Onderzoek Standaardbepaling TGN Fase 2

3.1.1 Training en selectie van beoordelaars

Voor het onderzoek van Fase 2 werden verspreid over 3 dagdelen in totaal 55 kandidaat-beoordelaars getraind. De training was vergeleken met de training van beoordelaars in Fase 1, zie paragraaf 2.1.1, uitgebreider en richtte zich zowel op het A1min- en A1-niveau als het A2-niveau. Ook in Fase 2 werd opnieuw gebruik gemaakt van videofragmenten. Daarnaast werd gewerkt met geluidsopnames van intakegesprekken. Het videomateriaal laat gesprekken zien tussen een moedertaalspreker van het Nederlands en inburgeraars (NT2 sprekers). Er werd gebruik gemaakt van videomateriaal van diverse gesprekken waarbij kandidaten met verschillende niveaus aan bod komen. De geluidsfragmenten bestonden uit een intake gesprek tussen de Nederlandstalige intaker en de NT2 leerders. De intakegesprekken werden gemaakt in het kader van het plaatsen van inburgeraars op inburgeringstrajecten.

Na de training hebben 3 kandidaat-beoordelaars zich teruggetrokken omdat ze de opdracht te moeilijk vonden, de opdracht voor hen technisch niet haalbaar was met de thuiscomputer of omdat ze niet aan de planning konden voldoen. In totaal hebben 52 kandidaat-beoordelaars het assessment afgelegd, hebben zich na dit assessment 7 kandidaat-beoordelaars teruggetrokken en hebben 6 kandidaat-beoordelaars een tweede assessment gemaakt. Uiteindelijk zijn 40 beoordelaars bij de beoordeling van kandidaten ingezet.

Een (eerste) groep van 20 beoordelaars beoordeelde 90 kandidaten die een TGN score rond de cesuur van het A1min-niveau hadden en 40 kandidaten die een TGN score rond de cesuur van het A1-niveau hadden. Een (tweede) groep van 20 beoordelaars beoordeelde 40 kandidaten die een TGN score rond de cesuur van het A1-niveau hadden en 80 kandidaten die een TGN score rond de cesuur van het A2-niveau hadden.

Aan voornoemde aantallen beoordelaars en kandidaten lagen de volgende overwegingen ten grondslag. Een aantal van 15 tot 30 beoordelaars wordt in de vakliteratuur over standaardbepaling als een voldoende aantal aangemerkt (Hambleton & Pitoniak, Educational Measurement, 2006, p. 452). Het aantal van 120/130 kandidaten dat een beoordelaar in dit onderzoek moest beoordelen is het maximum aantal kandidaten dat een beoordelaar in het beschikbare tijdsbestek redelijkerwijs geacht werd te kunnen beoordelen.

3.1.2 Beoordeling van kandidaten

De beoordelaars kregen het examen van de kandidaat-inburgeraar te horen. De prompts kregen de beoordelaars niet te horen maar ze kregen de prompts echter wel te zien om de inhoudelijke correctheid van de responsen op adequate wijze mee te kunnen nemen. De beoordelingstaak voor de beoordelaars bestond uit het nemen van de (dichotome) beslissing of een kandidaat wel of niet voldeed aan het beschreven niveau.

3.1.3 Selectie kandidaten onderzoek naar standaard A1min-niveau

Voor het onderzoek naar de standaard voor het A1min-niveau werd uit het bestand van de kandidaten die de TGN afgelegd hadden, een gestratificeerde random steekproef getrokken van kandidaten die op de TGN respectievelijk de (door het consortium oorspronkelijk toegekende) scores 16, 18, 20, 22, 24, 26, 28, 30 en 32 behaald hadden. Deze scorepunten of scoreniveaus zijn evenredig verspreid over het deel van de TGN schaal waarbinnen naar verwachting de standaard voor het A1min-niveau ligt. Voor deze 9 scorepunten werden in totaal 180 kandidaten, 20 kandidaten per scorepunt, geselecteerd. De selectie van 180 kandidaten was om de invloed van aberrante kandidaten, dat wil zeggen kandidaten die volgens een groot aantal beoordelaars beter of slechter zijn dan zij volgens hun TGN score zijn, op de uitkomsten van de standaardbepaling zoveel mogelijk te beperken. Aangezien een beoordelaar echter maximaal 90 kandidaten zou beoordelen, kon niet elke beoordelaar alle geselecteerde 180 kandidaten beoordelen. De restrictie van maximaal 90 beoordelingen per beoordelaar en een andere restrictie zoals 10 beoordelingen per scorepunt, vereiste de inzet van een incompleet gebalanceerd afnamedesign. Bijlage 2.1 bevat het afnamedesign voor de standaardbepaling van het A1min-niveau. Het design is incompleet omdat niet alle beoordelaars alle examens beoordelen waardoor niet in alle cellen van de matrix een dichotome beoordeling vermeld staat. Het afnamedesign is gebalanceerd waarmee bedoeld wordt dat alle beoordelaars een gelijk aantal van 90 beoordelingen geven (zie de laatste rij van bijlage 2.1) en dat elk scorepunt een (nagenoeg altijd) gelijk aantal van 10 beoordelingen heeft (zie de kolom met 'oordelen' in bijlage 2.1).

3.1.4 Selectie kandidaten onderzoek naar standaard A1-niveau

Voor het onderzoek naar de standaard voor het A1-niveau werden 160 kandidaten geselecteerd die op de TGN de (door het consortium oorspronkelijk

toegekende) scores 26, 28, 30, 32, 34, 36, 38 en 40 behaald hadden. Deze scorepunten zijn evenredig verspreid over het deel van de TGN schaal waarbinnen naar verwachting de standaard voor het A1-niveau ligt. Voor deze 8 scorepunten werden in totaal 160 kandidaten, 20 kandidaten per scorepunt, geselecteerd. De selectie van 160 kandidaten vond op dezelfde wijze en op basis van dezelfde overwegingen plaats als die van de selectie van kandidaten voor het onderzoek naar de standaard voor het A1min-niveau.

In verband met de omvang van de beoordelingstaak voor de beoordelaars, werden bij de beoordeling van kandidaten 40 beoordelaars ingezet. Van die 40 beoordelaars beoordeelden 20 beoordelaars 40 kandidaten. In bijlage 2.2.1 dat het ene deel van het incompleet gebalanceerd afnamedesign bevat, hebben bedoelde beoordelaars de nummers 1 t/m 20 en de kandidaten de nummers 1 t/m 80. In de cellen van de matrijs staan de dichotome beoordelingen van de beoordelaars vermeld.

De 20 andere beoordelaars beoordeelden ook 40 kandidaten. In bijlage 2.2.2 dat het andere deel van het incompleet gebalanceerd afnamedesign bevat, hebben de beoordelaars de nummers 21 t/m 40 en de kandidaten de nummers 81 t/m 160. In de cellen van de matrijs staan de dichotome beoordelingen van de beoordelaars vermeld.

3.1.5 Selectie kandidaten onderzoek naar standaard A2-niveau

Voor het onderzoek naar de standaard voor het A2-niveau werden kandidaten geselecteerd die op de TGN de (door het consortium oorspronkelijk toegekende) scores 34, 36, 38, 40, 42, 44, 46 en 48 behaald hebben. Deze scorepunten zijn evenredig verspreid over het deel van de TGN schaal waarbinnen naar verwachting de standaard voor het A2-niveau ligt. Voor deze 8 scorepunten werden in totaal 160 kandidaten, 20 kandidaten per scorepunt, geselecteerd. De selectie van 160 kandidaten vond op basis van dezelfde overwegingen plaats als die van de selectie van kandidaten voor het onderzoek naar de standaard voor het A1min-niveau.

Bijlage 2.3 bevat het incompleet gebalanceerd afnamedesign voor de standaardbepaling van het A2-niveau waarbij in de cellen van de matrijs de dichotome beoordelingen van de beoordelaars vermeld staan.

3.2 Resultaten Onderzoek Standaardbepaling TGN Fase 2

3.2.1 Resultaten Standaardbepaling A1min-niveau

Bijlage 2.1 bevat de beoordelingen door 20 beoordelaars van de examens van 180 kandidaten. In bijlage 2.1 kan men zien dat er bij alle scorepunten sprake is van afwijkende beoordelingen van kandidaten met hetzelfde aantal scorepunten. Bij scorepunt 16 kan men zien dat kandidaat 5 door 8 van de 10 beoordelaars als A1min-niveau of hoger beoordeeld wordt. De andere 19 kandidaten met een TGN score van 16 worden niet als A1min-niveau of hoger

beoordeeld. Vanaf scorepunt 26 neemt het aantal afwijkende beoordelingen enigszins toe. Verschillende kandidaten met dezelfde scorepunten worden door alle 10 beoordelaars zowel als niet van A1min-niveau of hoger en als wel van A1min-niveau of hoger beoordeeld.

In bijlage 2.1 kan men zien dat er per scorepunt in totaal 200 examens dichotoom beoordeeld zijn. Van deze 200 examens is het percentage examens berekend dat als A1min-niveau of hoger beoordeeld wordt. Tabel 3.1 bevat per scorepunt het percentage examens dat als A1min-niveau of hoger beoordeeld wordt.

Tabel 3.1: A1min-niveau

Score	Percentage
16	16
18	17
20	18
22	20
24	29
26	38
28	50
30	46
32	54

In Tabel 3.1 kan men lezen dat maximaal 38% van de examens van kandidaten met scorepunten 16 t/m 26 als A1min-niveau of hoger beoordeeld worden. Van de examens van kandidaten met scorepunt 28 wordt 50% als A1min-niveau of hoger beoordeeld. Van de examens van kandidaten met scorepunt 30 wordt 46% en van kandidaten met scorepunt 32 wordt 54% als A1min-niveau of hoger beoordeeld.

3.2.2 Resultaten Standaardbepaling A1-niveau

Bijlagen 2.2.1 en 2.2.2 bevatten de beoordelingen door 40 (20+20) beoordelaars van de examens van 160 kandidaten. Uit bijlage 2.2.1 en 2.2.2 blijkt dat er bij alle scorepunten sprake is van discrepanties tussen de beoordelingen van de kandidaten. Deze discrepanties zijn minder groot bij de lagere scorepunten, bijvoorbeeld scorepunt 26, dan bij de hogere scorepunten zoals scorepunt 40.

In bijlage 2.2.1 en bijlage 2.2.2 kan men lezen dat er per scorepunt in totaal 200 examens dichotoom beoordeeld zijn. Van deze 200 examens is het percentage examens berekend dat als A1-niveau of hoger beoordeeld wordt. Tabel 3.2 bevat per scorepunt het percentage examens dat als A1-niveau of hoger beoordeeld wordt.

Tabel 3.2: A1-niveau

Score	Percentage
26	18
28	23
30	27
32	23
34	31
36	23
38	32
40	32

In Tabel 3.2 kan men lezen dat 18% van de examens van kandidaten met scorepunt 26 als A1-niveau of hoger beoordeeld worden en dat dit laatste voor 32% van de examens van kandidaten met scorepunt 38 en 40 geldt.

3.2.3 Resultaten Standaardbepaling A2-niveau

Bijlage 2.3 bevat de beoordelingen door 20 beoordelaars van de examens van 160 kandidaten. In bijlage 2.3 kan men zien dat er bij alle scorepunten sprake is van afwijkende beoordelingen van kandidaten met hetzelfde aantal scorepunten. In het algemeen zijn de afwijkingen bij de lagere scorepunten wat groter dan bij de hogere scorepunten.

In bijlage 2.3 kan men lezen dat er per scorepunt in totaal 200 examens dichotoom beoordeeld zijn. Van deze 200 examens is het percentage examens berekend dat als A1min-niveau of hoger beoordeeld wordt. Tabel 3.3 bevat per scorepunt het percentage examens dat als A2-niveau of hoger beoordeeld wordt.

Tabel 3.3: A2-niveau

Score	Percentage
34	08
36	13
38	15
40	22
42	18
44	14
46	20
48	43

In Tabel 3.3 kan men lezen dat maximaal 22% van de examens van kandidaten met scorepunten 34 t/m 46 als A2-niveau of hoger beoordeeld worden. Van de examens van kandidaten met scorepunt 48 wordt 43% als A2-niveau of hoger beoordeeld.

4. Conclusies en Aanbevelingen

4.1 Conclusies

In paragraaf 3.2 zijn de resultaten van het onderzoek naar de standaardbepaling voor het A1min-niveau, A1-niveau en het A2-niveau beschreven. Op basis van die resultaten kunnen de volgende conclusies getrokken worden.

Voor wat betreft het A1min-niveau kunnen op basis van de resultaten van het onderzoek twee conclusies getrokken worden:

1. De resultaten bevestigen de conclusie van eerdere onderzoeken van TNO en het RCEC dat de oorspronkelijke cesuur van 16 scorepunten op de TGN voor het A1min-niveau te soepel was.
2. De resultaten vermeld in Tabel 3.1 maken echter niet overtuigend duidelijk wat die cesuur wel zou moeten zijn. Vandaar dat zowel door de onderzoekers van het RCEC als TNO aanvullende statistische analyses met behulp van zogenaamde meerniveau logistische regressiemodellen uitgevoerd zijn. Voor een beschrijving van het gehanteerde model en de analyse wordt verwezen naar Bijlage 3. Deze analyses leveren schattingen van de gemiddelde score of cesuur op en schattingen van standaardfouten met behulp waarvan betrouwbaarheidsintervallen bepaald kunnen worden. Afhankelijk van de foutenbronnen waarmee rekening gehouden wordt, resulteren de analyses in verschillende betrouwbaarheidsintervallen en dus ook in verschillende ondergrenzen. De resultaten van de analyses laten zien dat een ondergrens van 25 à 26 scorepunten het meest aannemelijk is. Als met alle foutenbronnen rekening gehouden wordt, is een ondergrens van 21 à 22 scorepunten ook mogelijk.

Voor wat betreft het A1-niveau maken de resultaten vermeld in Tabel 3.2 niet overtuigend duidelijk welke score kandidaten op de TGN zouden moeten behalen om als A1-niveau of hoger beoordeeld te worden. Verwacht had mogen worden dat een toename van het aantal scorepunten gepaard gegaan zou zijn met een toename van de percentages maar dat blijkt in de gevonden data nauwelijks het geval te zijn. Gegeven de range van scorepunten lijkt een cesuur in de buurt van de hogere scorepunten, 38 en 40, echter aannemelijker dan een cesuur in de buurt van de lagere scorepunten, 26 en 28.

Voor wat betreft het A2-niveau kan men in Tabel 3.3 lezen dat maximaal 22% van de examens van kandidaten met scorepunten 34 t/m 46 als A2-niveau of hoger beoordeeld worden en dat 43% van de examens van kandidaten met scorepunt 48 als A2-niveau of hoger beoordeeld worden. Ook voor het A2-niveau geldt dat de resultaten niet overtuigend duidelijk maken welke score kandidaten op de TGN zouden moeten behalen om als A2-niveau of hoger beoordeeld te worden. Vanwege het relatief hoge percentage bij scorepunt 48 lijkt een cesuur in de buurt van 48 scorepunten aannemelijk te zijn.

4.2 Aanbevelingen

Naar aanleiding van eerder onderzoek van TNO en het RCEC is de oorspronkelijke cesuur van het A1-min niveau met ingang van 15 maart 2008 verhoogd tot 22 scorepunten en zijn de andere niveaus overeenkomstig aangepast. Aanbevolen wordt om 22 scorepunten als cesuur voor het A1min-niveau te handhaven. Hiervoor zijn twee redenen aan te geven. In de eerste plaats de resultaten van het onderhavige onderzoek die aangeven dat een cesuur van 22 scorepunten, de ondergrens van het meest conservatieve betrouwbaarheidsinterval, een voorzichtige keuze is die methodologisch te verantwoorden valt. In de tweede plaats valt een cesuur van 22 scorepunten ook inhoudelijk te verantwoorden, dat wil zeggen dat een cesuur van 22 scorepunten door inhoudelijk deskundigen als bovengrens beschouwd wordt van de productieve en receptieve woordenschat die van kandidaten van dit niveau verwacht mag worden. In termen van de opdrachten en vragen van de TGN betekent een cesuur van 22 scorepunten dat kandidaten gemiddeld 26% van de herhaalopdrachten en 25% van de woordvragen goed moeten beantwoorden. Een hoger beheersingspercentage past niet meer in de omschrijving van niveau A1min zoals opgesteld door de Adviescommissie Normering Inburgeringseisen en zoals die is afgestemd op de overige niveaus van de CEF-schaal.

Om de continuïteit van de scoreschaal niet te verstoren, wordt aanbevolen de cesuren voor de andere niveaus te handhaven op de niveaus zoals die met ingang van 15 maart 2008 gehanteerd werden.

Bijlage 1 Onderzoek Standaardbepaling TGN Fase 1

Bijlage 1.1 Data TGN 2008 Fase 1, A1Min-niveau

knd-	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	Tot	
Beo	11	11	11	14	14	14	17	17	17	20	20	20	23	23	23	26	26	26	29	29	29	29	Tot
A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1
min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min	min
1	7	8	16	20	2	17	9	19	10	3	15	11	21	4	12	5	14	6	18	13	13	13	
1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	1	1	1	0	0	1	5	
2	0	0	0	1	0	0	1	0	0	1	1	0	0	1	1	1	1	1	1	0	1	10	
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	
4	0	0	0	1	0	0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	4	
5	0	0	0	1	0	0	1	0	0	0	0	0	0	1	0	1	1	1	0	0	0	5	
6	0	0	0	0	0	0	0	0	0	1	0	0	0	1	1	1	1	1	1	0	1	7	
7	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	1	1	0	5	
8	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	2	
9	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	1	1	1	6	
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
11	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	3	
12	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	1	1	1	1	0	1	5	
13	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	1	1	1	1	0	1	6	
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
15	0	0	0	1	0	0	0	1	0	1	1	0	0	1	1	1	1	1	1	0	1	11	
16	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	3	
17	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	3	
18	1	1	0	1	0	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	1	17	
19	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0	1	1	0	0	0	0	3	
20	0	0	0	1	0	0	1	0	0	0	1	0	0	1	0	1	1	1	1	1	1	10	
21	0	0	0	1	0	0	0	0	1	0	1	0	1	1	1	1	1	1	1	1	1	12	
22	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0	1	0	0	0	0	0	3	
23	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	1	1	0	1	1	6	
24	0	1	0	0	0	0	0	0	0	1	0	0	0	1	0	1	0	0	0	0	0	4	
25	0	0	0	0	0	0	1	0	0	0	0	0	0	1	1	1	1	1	0	1	0	6	
Tot	1	2	0	7	0	1	5	2	2	1	10	7	1	3	21	8	23	18	9	5	11	5,48	

Proportie 0,04

0,11

0,12

0,24

0,33

0,65

0,33

Bijlage 1.2 Data TGN 2008 Fase 1, A1-niveau

knd-	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	som																							
Beo	20	20	20	23	23	23	26	26	26	29	29	29	32	32	32	35	35	35	38	38	38	38																							
7	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	A1	9																						
12	1	10	2	18	15	11	3	6	14	19	8	16	21	4	17	20	13	5	9	7	5	0	6	5	7	4	6	9	4	4	3	8	1	4	3	9	6	1	8	14	7	7	0	6	5,36
1	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	1	1	1	7																							
2	0	0	0	0	0	1	0	0	0	0	0	1	0	0	1	0	1	0	0	1	1	5																							
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0																							
4	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	1	1	1	0	6																							
5	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	0	1	1	1	5																							
6	0	0	0	1	0	0	0	0	0	0	0	1	1	0	0	1	0	1	1	1	0	7																							
7	0	0	0	0	0	0	0	0	0	1	0	0	1	0	0	1	0	1	0	0	0	4																							
8	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	1	0	1	1	1	1	6																							
9	0	0	0	0	1	1	0	0	0	1	0	1	0	1	0	1	0	1	1	1	1	9																							
10	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	1	4																							
11	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0	1	1	4																							
12	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	1	3																							
13	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	1	1	1	1	1	1	8																							
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1																							
15	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	0	0	1	0	4																							
16	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	1	0	3																							
17	0	0	0	0	1	1	0	0	0	0	0	1	0	1	0	1	0	1	1	1	1	9																							
18	0	0	0	1	1	0	0	0	0	0	0	0	0	0	1	0	0	1	1	1	1	6																							
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1																							
20	0	0	0	0	1	0	0	0	1	0	0	1	0	0	1	0	1	0	1	1	1	8																							
21	0	0	0	0	1	1	0	0	1	0	1	1	1	1	1	1	1	1	1	1	1	14																							
22	0	0	0	0	1	0	0	0	0	0	0	1	0	1	0	1	0	1	1	1	1	7																							
23	0	0	0	0	0	0	0	0	1	0	0	1	0	1	0	1	0	1	1	1	1	7																							
24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0																							
25	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	1	0	1	1	1	1	6																							
Tot	0	0	1	3	11	5	0	0	4	1	2	14	2	11	7	4	17	2	15	20	15	5,36																							

Proportie 0,01

0,25

0,05

0,23

0,27

0,31

0,67

Bijlage 2 Onderzoek Standaardbepaling TGN Fase 2

num. Pie	audiofile A2	beo 21	beo 22	beo 23	beo 24	beo 25	beo 26	beo 27	beo 28	beo 29	beo 30	beo 31	beo 32	beo 33	beo 34	beo 35	beo 36	beo 37	beo 38	beo 39	beo 40	# oordelen	score		Proportie	
91	score 42	0	0	0	0	0	0			0	0	0		0	0	0	0					10	0			
92	score 42	0	0	0			0			0	0			0	0	0	0			0	0	10	0			
93	score 42			0				0	0	0	0	0	0	0	0	0	0	0				10	0			
94	score 42			0	0	0	0	0			0	0					0	0				10	0			
95	score 42		0			0	0	0				0	0			0	0	0	0			10	0			
96	score 42	0	0	0	0	0	0				0	0			1	0			1	0		10	2			
97	score 42	0	0	0	0	0				0			0			0	0	0		0	0	10	0			
98	score 42	0	0	0			0			0	0	0	0	0	0				0	0	0	10	0			
99	score 42	0	0		0			0	0	0		0	0						0	0	0	10	0			
100	score 42	0		1	0					0	1		0	0				0	1		1	10	4			
101	score 44	0								0				0	0	0	0	0			0	10	0	score 44	0,14	
102	score 44		0	0		0			0		0	0	0	0			0					10	0			
103	score 44	0	0	0					0	0	0	0	0				0					10	0			
104	score 44		0	0		0			0	0				0				0	0	0		10	0			
105	score 44	0			0				0	0				0		0	0	0	1		0	10	1			
106	score 44			0				0			0	0	0	0		0		0	0	0	0	10	0			
107	score 44	0		1	0		0	1	0	1			0	1			1					10	5			
108	score 44	0	0		0	0	0			1							0	0	0	1	0	10	2			
109	score 44		0		0	0	0	0			0					0	0	0				10	0			
110	score 44	0			0	0	0			0	0	0				0		0			0	10	0			
111	score 44	0		0	0	0			0		0	0	0			0						10	0			
112	score 44	0	0							0		1		0	1	0	0			1		10	3			
113	score 44			0	0		0	0		0			0	1		0	0				0	10	1			
114	score 44			0			1			0	1		0	1	0	0			1	1		10	5			
115	score 44			0	0	0	0	0				0	0		0	0			0	0		10	0			
116	score 44		0		0				0		0	0	0					0		0	0	10	0			
117	score 44	0	0	0		1	0			0		1						1			0	10	3			
118	score 44		0			1	1	0				1	0	0	1		0	1				10	5			
119	score 44		0		0	0	0	0		0	0				0	0		0	0			10	0			
120	score 44	0			0	0			0	0	1			1	0	0				1		10	3			
121	score 46	0			1		1	1			1	0	1				1	1	1			10	8	score 46	0,2	
122	score 46			1	1	1			0		1	0					0	1		1	1	10	7			
123	score 46		0	0					0	0	0	0				0	0	0	0			10	0			
124	score 46	0		0					0	0	0	0				0		1		0		10	1			
125	score 46	0			0	0	0			0		0	1	1	0	0						10	2			
126	score 46		0	0				0	0	0	0									0	0	10	0			
127	score 46	0	0	0			0	0		0		0			1	0	0	1				10	2			
128	score 46					1	1	0	1			0			1	1	1		1	1		10	8			
129	score 46	0	0			0	0	0		0	0	0					0					10	0			
130	score 46				0	0	0	0		0	0						0	0				10	0			
131	score 46				0	0	0		0	0		0	0				0					10	0			
132	score 46	0		0	0	0							1	0				0	0	0	0	10	1			
133	score 46	0		0	0							0				0	0	0	0	1	0	10	1			
134	score 46	0	0		0	0				0				0	0	0		0	0			10	0			
135	score 46		0	0	0	1			0		0			0	0	0						10	1			
136	score 46		0						0	0	0	0	0	0	0				0			10	0			
137	score 46	0	0	1	0			0	0	0	0	0		1	0							10	2			
138	score 46	0	1				0	0	0			0	0	1	0							10	2			
139	score 46		1	0	0	0		1		1	1					0				1	0	10	5			
140	score 46		0	0				0	0		0	0						0	0	0		10	0			
141	score 48	0					0			0	0	0	0			0	0				0	10	0	score 48	0,437186	
142	score 48	0	0	0		0				0			0		1	0	0	0				10	1			
143	score 48	0	0		0	0			0	0			0	1		0		0				10	1			
144	score 48	0						0			0	0		1	1			0	1	1		0	10	4		
145	score 48		1		0		0			0		1					1	1	1	1		10	7			
146	score 48				0	0	0			1					1	0	0	1	1	1		10	5			
147	score 48			0	0	0	0				0			0	0	0		0			0	10	0			
148	score 48			1	0	0	0			0	0		0	1			1			1		10	4			
149	score 48			0	0	0			0		0	0	1	1				0				10	2			
150	score 48	0		1			1	1	0	1	1	1					1		1			10	7			
151	score 48	0			1		0	1		1	1	0			1	0					1	10	6			
152	score 48		0					0		0	0	0	0	0	1	0					1	10	2			
153	score 48		1	1			1				0	1	1	1		1		1	1		0	10	9			
154	score 48	0		0	0				0		0	0		0		0		0	0			10	0			
155	score 48		0				0		0	0	0			1	0			0		0	0	10	1			
156	score 48		0		1	1				0	1	0		1					1	1	0	10	6			
157	score 48	1	1	1	1		1	1	0				1						1	1		10	9			
158	score 48	0	1	1		1		1	0	0				1			0	1				10	6			
159	score 48				1					1	1	1	0	1						1	1	9	8			
160	score 48	0	1	1	1	1		1		1								1	1	1	1	10	9			

Bijlage 3 Betrouwbaarheidsinterval standaard A1min-niveau

Voor de bepaling van de standaard voor het A1min-niveau is met behulp van het computerprogramma MLWiN een meerniveau logistisch regressiemodel geschat. In het gehanteerde model is een modellering toegepast waarbij de beoordelaars (en hun oordelen over een kandidaat) genest zijn binnen de kandidaat. In het model zijn de beoordelaars niveau 1 en de kandidaten niveau 2. Als voorspeller op niveau 2 wordt de score van de kandidaat gebruikt. De standaardfouten van de modelparameters zijn bepaald met een parametrische bootstrap procedure.

In termen van MLWiN is het regressiemodeld geformuleerd als

$$C4_{ij} \sim \text{Binomial}(\text{const}_{ij}, \pi_{ij})$$

$$\text{logit}(\pi_{ij}) = \beta_{0j} \text{const} + 0.131(\sim 0.013) \text{score}_j$$

$$\beta_{0j} = -3.954(\sim 0.319) + u_{0j}$$

$$\begin{bmatrix} u_{0j} \end{bmatrix} \sim N(0, \Omega_u) : \Omega_u = \begin{bmatrix} 0.023(\sim 0.030) \end{bmatrix}$$

$$\text{var}(C4_{ij} | \pi_{ij}) = \pi_{ij}(1 - \pi_{ij}) / \text{const}_{ij}$$

Indien rekening gehouden wordt met de onzekerheid van de verschillende modelparameters, geeft onderstaande tabel schattingen voor de gemiddelde score en drie ondergrenzen van 95% betrouwbaarheidsintervallen.

grens CI normaal		1,96
	schatting	SE
b0	-3,954	0,319
b1	0,131	0,013
	0,023	0,03
gemiddelde score	30,18321	
ondergrens totaal	20,98466	
ondergrens obv b0 +beoord.	25,06626	
ondergrens obv b0	25,41038	

Wat betreft de ondergrenzen van de betrouwbaarheidsintervallen laat bovenstaande tabel zien dat wanneer alleen rekening wordt gehouden met de

foutenbron op basis van b_0 , de ondergrens van het betrouwbaarheidsinterval gelijk is aan 25.4, en dat de ondergrens op basis van b_0 en de verdeling van de beoordelaars gelijk is aan 25.1. De 'ondergrens totaal' in de tabel bevat alle bronnen van onzekerheid waarbij elke foutenbron een afwijking tot de puntschatting heeft van 1.96 standaarddeviatie. Dit laatste resulteert in een ondergrens van het betrouwbaarheidsinterval die gelijk is aan 21.0.