

# Wetenschappelijke verantwoording Spelling 3.0 digitaal voor groep 7

Aanvulling bij de wetenschappelijke verantwoording van de LVS-toetsen  
Spelling 3.0 voor groep 7

Marieke Tomesen, Jasper Wouda en Linda Horsels

[cito.nl](https://cito.nl)





# **Wetenschappelijke verantwoording Spelling 3.0 digitaal voor groep 7**

Aanvulling bij de wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7

Cito | Primair en speciaal onderwijs

Marieke Tomesen  
Jasper Wouda  
Linda Horsels

© Cito B.V. Arnhem (2019)

Niets uit dit werk mag zonder voorafgaande schriftelijke toestemming van Cito worden openbaar gemaakt en/of verveelvoudigd door middel van druk, fotokopie, scanning, computersoftware of andere elektronische verveelvoudiging of openbaarmaking, microfilm, geluidskopie, film- of videokopie of op welke wijze dan ook. welke wijze dan ook.

# Inhoud

<b>1</b>	<b>Inleiding</b>	<b>5</b>
<b>2</b>	<b>Uitgangspunten van de toetsconstructie</b>	<b>7</b>
<b>3</b>	<b>Beschrijving van de toetsen</b>	<b>9</b>
3.1	Opbouw en structuur van de toetsen	9
3.2	Inhoudsverantwoording	10
3.2.1	Domeinbeschrijving en uitwerking in spellingcategorieën	10
3.2.2	Itemconstructie, onderzoeken en selectie van opgaven	10
3.3	Statistische beschrijving	15
<b>4</b>	<b>Kalibratie en normering</b>	<b>19</b>
4.1	Rationale van de kalibratieonderzoeken	19
4.2	Kalibratieonderzoek digitale items	19
4.2.1	Opzet van het kalibratieonderzoek voor de digitale items	19
4.2.2	De stappen in de kalibratie	23
4.2.3	Toetsing van het IRT-model	24
4.2.4	Totale kalibratie per groep	25
4.3	De normering	29
4.3.1	Opzet	30
4.3.2	Representativiteit	31
4.3.3	Normeringsresultaten	32
<b>5</b>	<b>Betrouwbaarheid en meetnauwkeurigheid</b>	<b>33</b>
5.1	Betrouwbaarheid	33
5.2	Nauwkeurigheid	34
<b>6</b>	<b>Validiteit</b>	<b>41</b>
<b>7</b>	<b>Samenvatting</b>	<b>43</b>
<b>8</b>	<b>Aanvullende literatuur</b>	<b>45</b>
<b>Bijlagen 47</b>		
1	Moeilijkheid van opgaven per taak in Spelling 3.0 digitaal groep 7	48
2	Klassieke en IRT-indices van de opgaven in digitale toetsen Spelling 3.0 groep 7	52



# 1 Inleiding

Deze *Aanvulling bij de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7* heeft uitsluitend betrekking op de *digitale* toetsen Spelling 3.0 voor groep 7.

De inhoud van deze digitale toetsen komt grotendeels overeen met de inhoud van de papieren toetsen voor groep 7. Vandaar dat we voor de inhoudelijke aspecten grotendeels verwijzen naar de oorspronkelijke Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018). Op punten waar de digitale toetsen afwijken van de papieren toetsen, gaan we in deze aanvulling in.

Tezamen met de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018) en de inhoud van het (digitale) toetspakket Spelling groep 7 (Cito, 2017) levert deze aanvulling alle informatie die nodig is voor een snelle en efficiënte beoordeling van de kwaliteit van de digitale toetsen Spelling 3.0 groep 7. Het genoemde materiaal maakt een beoordeling van de digitale toetsen Spelling groep 7 mogelijk op de volgende aspecten:

- Uitgangspunten van de toetsconstructie
- Uitgangspunten van de toetsconstructie
- Kwaliteit van het toetsmateriaal
- Kwaliteit van de handleiding
- Normen
- Betrouwbaarheid
- Validiteit

Het laatstgenoemde aspect betreft alleen begripsvaliditeit en geen criteriumvaliditeit. Omdat de toetsen van het LVS niet bedoeld zijn voor 'voorspellend gebruik' is criteriumvaliditeit niet van toepassing.

Deze aanvulling heeft met name betrekking op de normen (hoofdstuk 4) en de betrouwbaarheid en meetnauwkeurigheid (hoofdstuk 5). Voor de uitgangspunten van de toetsconstructie (hoofdstuk 2 en 3) en de begripsvaliditeit (hoofdstuk 6) van de toetsen Spelling 3.0 verwijzen we naar de oorspronkelijke Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018). De kwaliteit van het toetsmateriaal en van de handleiding is te bepalen door kennis te nemen van de inhoud van de (digitale) toetspakketten.





## 2 Uitgangspunten van de toetsconstructie

Voor de uitgangspunten van de toetsconstructie verwijzen we naar de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018). Alles wat in hoofdstuk 2 van de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 wordt gezegd over de meetpretentie, het gebruiksdoel en de functie van de toetsen, is ook van toepassing op de digitale toetsen.

Meetpretentie, gebruiksdoel en functie zijn identiek voor de papieren en digitale toetsen. De papieren en digitale toetsen zijn immers op dezelfde manier opgebouwd en in grote lijnen gelijk aan elkaar. Voor zowel de papieren als de digitale toetsen geldt het volgende:

- ze meten de actieve spelling doordat de leerling woorden moet opschrijven cq. intypen;
- ze zijn bestemd voor leerlingen in groep 7 van het basisonderwijs en voor leerlingen in het speciaal basisonderwijs en in het speciaal onderwijs cluster 1, cluster 2 (leerlingen met een taalontwikkelingsstoornis) en cluster 4;
- voor zowel ‘midden leerjaar’ (half januari/half februari) als voor ‘einde leerjaar’ (juni) zijn populatieparameters bepaald;
- ze kunnen ook gebruikt worden voor leerlingen in andere leerjaren die werken op het niveau van groep 7;
- de toetsen zijn niet geschikt voor leerlingen met een (tijdelijk) beperkt gehoor;
- de toetsen hebben twee doelen: niveaubepaling en progressiebepaling;
- de gemaakte fouten kunnen geanalyseerd worden met het oog op het aanbieden van gerichte remediëring.

Van alle papieren toetsen Spelling 3.0 zijn ook digitale varianten beschikbaar. Dit betekent dat er voor groep 7 een digitale toets M7 niet-werkwoorden, een digitale toets E7 niet-werkwoorden, een digitale toets M7 werkwoorden en een digitale toets E7 werkwoorden beschikbaar is. De papieren en digitale toetsen van een bepaald afnamemoment zijn uitwisselbaar. Dat betekent dat de leerkracht op een afnamemoment zelf kan kiezen of hij/zij de leerling een papieren of digitale toets laat maken. De keuze heeft geen invloed op het volgen van de ontwikkeling van de spellingvaardigheid. De scores op de papieren en digitale toetsen zijn namelijk uitwisselbaar. Ze zijn echter niet identiek, wat betekent dat eenzelfde aantal goed tot een andere vaardigheidsscore leidt. De omzetting van vaardigheidsscore naar niveau is echter hetzelfde.

De theoretische inkadering van de toetsen - zowel inhoudelijk als psychometrisch (zie paragraaf 2.4 van de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018) - geldt zowel voor de papieren toetsen als de digitale toetsen.



## 3 Beschrijving van de toetsen

### 3.1 Opbouw en structuur van de toetsen

Het digitale toetspakket Spelling 3.0 voor groep 7 uit het Cito Volsysteem primair en speciaal onderwijs bevat – net als de papieren toetsen voor groep 7 – in totaal vier toetsen: twee toetsen Spelling niet-werkwoorden (M7 en E7) en twee toetsen Spelling werkwoorden (M7 en E7). Voor de toetsen M7 zijn de populatieparameters bepaald op ‘midden leerjaar’ groep 7, voor de toetsen E7 op ‘einde leerjaar’ groep 7.

De digitale varianten van de toetsen bevatten net als de papieren varianten twee taken van 25 opgaven, met uitzondering van de toets Spelling werkwoorden M7; deze bestaat uit twee taken van 20 opgaven. De papieren en digitale toetsen zijn op dezelfde manier opgebouwd. Het opgaventype (zinsdictee) is hetzelfde. Ook het toetsen op maat is evengoed mogelijk bij de digitale toetsen. We verwijzen daarom hier naar hoofdstuk 3 van de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018) voor een nadere toelichting.

De opgaven van de papieren en de digitale toetsen representeren dezelfde spellingcategorieën.

De dicteewoorden die in een digitale toets voorkomen kunnen echter verschillen van de papieren toets van hetzelfde niveau. De papieren en digitale toetsen zijn dus niet identiek, maar hebben wel een grote overlap.

De afname, scoring en verwerking van de resultaten verschilt van die van de papieren toetsen. We lichten ze hieronder toe.

#### *Afname*

De digitale toetsen worden individueel op de computer, laptop of chromebook gemaakt. Afhankelijk van het aantal beschikbare computers kunnen meerdere leerlingen gelijktijdig aan dezelfde toets werken.

De leerling krijgt voorafgaand aan de toets een uitleg over het maken van de digitale opgaven en maakt twee oefenopgaven. Op die manier raakt de leerling vertrouwd met het type opgave.

Bij het maken van de instructie is rekening gehouden met speciale leerlingen; zo is de instructie bijvoorbeeld kort en zijn samengestelde zinnen zo veel mogelijk vermeden. Dergelijke uitgangspunten gaan niet ten koste van de reguliere leerlingen.

Daarna maakt de leerling de toetsopgaven. Bij de toets Spelling niet-werkwoorden ziet de leerling op het scherm een antwoordvak. De computer leest de opgave voor. Daarna verschijnt de cursor in het antwoordvak, zodat de leerling het dicteewoord kan intikken. Bij de toets Spelling werkwoorden ziet de leerling het hele werkwoord, en daaronder de zin met een invulvak. Het is niet mogelijk om tijdens het afspelen van de audio al een antwoord in te tikken. De leerling kan, indien gewenst, de opgave nog een keer afspelen door op de ‘play-knop’ te klikken. Ook kan hij het volume aanpassen met de knop ‘luidspreker’; deze staat rechts boven in het scherm. Naast deze knop staat de knop voor het overzichtsscherm. Zo ziet de leerling in één oogopslag welke opgaven hij al gemaakt heeft, wat hij heeft ingevuld én welke opgaven hij nog moet maken. De leerling kan zijn antwoord aanpassen door de knop ‘backspace’ op het toetsenbord te gebruiken. Bij de toets E7 niet-werkwoorden staan rechts in beeld zeven knoppen met letters met een leesteken: è, é, ë, ï, ü, 's en s'. De leerling kan er ook voor kiezen om het toetsenbord te gebruiken om letters met een leesteken op te nemen. Bij de overige toetsen zijn deze knoppen niet van toepassing.

Als de leerling een antwoord heeft ingetypt, klikt hij op de knop ‘verder’. De gemaakte opgave krijgt dan een andere kleur (lichtgrijs). Het is mogelijk om via de navigatiebalk door de toets te navigeren en de opgaven in een andere volgorde te maken. Aan het einde van de toets kan de leerling zijn antwoorden controleren in het zogenaamde ‘overzichtsscherm’. Indien gewenst kan de leerling hier op een antwoord klikken en deze aanpassen. Pas als alle opgaven gemaakt zijn, kan hij de toets inleveren.

In de toetsmap Spelling 3.0 groep 7 (Cito, 2017) is een inhoudelijke handleiding opgenomen behorend bij de papieren en digitale toetsen. Hierin staan de uitgebreide afname-instructies voor de leerkracht. Daarnaast is er een technische digitale handleiding voor alle leerjaren (Cito, 2019), die voor scholen via Cito Portal gedownload kan worden.

Bij de digitale versies van de toetsen worden de antwoorden van de leerlingen door de computer gescoord en hoeft de leerkracht de toetsen dus niet zelf na te kijken. De leerkracht kan ervoor kiezen om een foutenanalyse uit te draaien waarin hij kan zien in welke spellingcategorieën de leerling veel fouten maakt.

Het maakt voor de resultaten niet uit of leerlingen de papieren of de digitale toetsen maken. De opgaven uit de papieren en de digitale toetsen liggen op één vaardigheidsschaal, waardoor de toetsresultaten onderling uitwisselbaar zijn. Bij de keuze voor de afname van ofwel de papieren ofwel de digitale toets kunnen verschillende overwegingen een rol spelen. Dit zijn overwegingen van zowel praktische aard (bijvoorbeeld de aanwezigheid van voldoende computers) als van meer inhoudelijke aard. Vooral voor leerlingen met concentratieproblemen, leerlingen die langzamer of juist veel sneller dan gemiddeld werken en leerlingen die afwezig waren bij de klassikale afname, kan een individuele, digitale afname prettig zijn. De leerling moet wel voldoende computervaardig zijn om woorden te kunnen intypen bij de digitale toets.

### *Scoring*

De digitale toetsen worden geautomatiseerd nagekeken. De toetsscore wordt automatisch omgezet naar de bijbehorende vaardigheidsscore met een score-interval ofwel betrouwbaarheidsinterval.

### *Verwerking resultaten*

Met het Computerprogramma LOVS kunnen allerlei rapportages, zoals leerlingrapporten en groepsoverzichten, en een foutenanalyse worden opgevraagd.

## **3.2 Inhoudsverantwoording**

### 3.2.1 Domeinbeschrijving en uitwerking in spellingcategorieën

Aan de digitale toetsen ligt dezelfde domeinbeschrijving en uitwerking in spellingcategorieën ten grondslag als aan de papieren toetsen. Hiervoor verwijzen we naar paragraaf 3.2.1 van de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018). De voor groep 7 afzonderlijk uitgewerkte overzichten van spellingcategorieën voor niet-werkwoorden cq. werkwoorden (zie tabel 3.1a cq. 3.1b) vormde de basis voor de itemconstructie en de selectie van items in de definitieve toetsen.

### 3.2.2 Itemconstructie, onderzoeken en selectie van opgaven

#### *Itemconstructie*

Er heeft geen speciale itemconstructie plaatsgevonden voor de digitale toetsen. Bij de digitale toetsen putten we uit de items van de itembank die is gevormd bij de constructie van de papieren toetsen. De papieren opgaven werden 'omgebouwd' tot een digitaal afneembare versie. De instructies en dicteeopgaven zijn ingesproken door een voice-over aan de hand van scripts. Een toetsdeskundige was aanwezig bij de opnames om de gesproken teksten direct te beoordelen en waar nodig bij te sturen. Voorafgaand aan de opnames bespraken zij aan de hand van voorbeeldaudio het spreektempo. Bij twijfel over de uitspraak van een dicteewoord werd het uitspraakwoordenboek geraadpleegd (zie [www.woorden.org](http://www.woorden.org)).

#### *Samenstelling definitieve toetsen*

In januari 2017 (niet-werkwoorden), juni 2017 (niet-werkwoorden en werkwoorden) en januari 2018 (werkwoorden) vonden kalibratieonderzoeken plaats voor de digitale items (zie hoofdstuk 4). Bij de analyses is de kwaliteit van de afzonderlijke items en de totale verzameling voor een afnamemoment in

kaart gebracht. Itemparameters en discriminatieparameters zijn geschat. Bij de analyses van de antwoorden van de leerlingen op de opgaven is nagegaan of de papieren items en de digitale items dezelfde vaardigheid meten en op dezelfde schaal passen. Dat bleek het geval te zijn.

Zie voor een uitgebreide verantwoording hoofdstuk 4. Voor de digitale items zijn eigen moeilijkheids- en discriminatieparameters geschat. Het is immers niet noodzakelijkerwijs zo dat de papieren versie en digitale versie van een item precies even moeilijk zijn en/of evengoed discrimineren. Met andere woorden: het is zeer wel mogelijk dat de papieren en digitale versie van hetzelfde item in deze kalibratie verschillende itemparameters krijgen toegekend.

Voor het samenstellen van de definitieve digitale toetsen zijn de volgende uitgangspunten gehanteerd:

- De digitale toetsen bevatten bij voorkeur precies dezelfde aantallen opgaven per categorie als de papieren toetsen.
- Van elke categorie zijn minimaal drie opgaven in een toets opgenomen.
- De digitale toetsen bevatten voor de meerderheid dezelfde items als de papieren toetsen van hetzelfde niveau.
- Een zelfde opgave mag maximaal twee keer voorkomen in het digitale volgsysteem Spelling 3.0.
- Er worden geen opgaven met DIF ten opzichte van de papieren toetsen opgenomen. Er is bij een digitale opgave sprake van DIF wanneer het bij gefixeerde itemparameters in de papieren versie niet mogelijk is het digitale item op dezelfde schaal te kalibreren.
- De gemiddelde moeilijkheid van de digitale toetsen is zoveel mogelijk gelijk aan die van de papieren toetsen vanwege eenzelfde toetsbeleving voor de leerlingen.
- De items van de digitale toetsen verwijzen naar precies dezelfde categorieën als de papieren toetsen.
- Net als bij de papieren toetsen komen in principe alle opgaven met een acceptabele moeilijkheid (in klassieke termen een p-waarde tussen 0,40 en 0,90) die door de betere spellers significant vaker goed worden gemaakt dan door de minder goede spellers (rir vanaf 0,20) in aanmerking voor opname in de definitieve digitale toetsen Spelling.

Aan de voorwaarden waarnaar deze uitgangspunten verwijzen, hebben we in grote lijnen kunnen voldoen. De digitale toetsen M7 en E7 Spelling niet-werkwoorden bevatten unieke opgaven en hebben grote overlap met de corresponderende papieren toetsen. De aantallen opgaven per categorie zijn bij de toetsen Spelling niet-werkwoorden in grote lijnen gelijk aan die van de papieren variant. Bij de toets M7 niet-werkwoorden zien we verschillen bij de categorieën 25, 32, 34 en 41, bij de toets E7 niet-werkwoorden alleen bij categorie 25 en 45. De afwijking bestaat elke keer uit één opgave meer of minder ten opzichte van de papieren toets. Er is wel vastgehouden aan het principe dat elke categorie met minimaal drie opgaven in de toets vóórkomt. De precieze aantallen opgaven per categorie staan in tabel 3.1a.

Tabel 3.1a Spellingcategoriegrenzen in de digitale toetsen Spelling 3.0 groep 7 niet-werkwoorden  
(met de aantallen in papieren toetsen tussen haakjes)

Spelling niet-werkwoorden			
Categorie	Omschrijving	M7	E7
20/21	woorden met open/gesloten lettergreep	4 (4)	4 (4)
24	woorden met -ig(-) en -lijk(-)	4 (4)	–
25	woorden waarin /ie/ geschreven wordt als i	4 (3)	4 (3)
26	woorden waarin /s/ geschreven wordt als c	–	4 (4)
27	woorden waarin /k/ geschreven wordt als c	–	4 (4)
28	woorden beginnend met 's of eindigend op 's	–	4 (4)
29	woorden met -tie(-) waarin t klinkt als (t)s	3 (3)	3 (3)
30	woorden met -heid of -teit	3 (3)	–
31	leenwoorden waarin /zj/ geschreven wordt als g(e)	4 (4)	–
32	leenwoorden waarin /sj/ geschreven wordt als ch (nieuw)	4 (3)	–
33	woorden met -b(-)	4 (4)	3 (3)
34	woorden met (-)y(-)	3 (4)	3 (3)
35	woorden met een trema	–	4 (4)
38	woorden met of zonder een hoofdletter	3 (3)	–
41	woorden waarin /t/ geschreven wordt als th	3 (4)	4 (4)
42	woorden met -isch(e)	4 (4)	4 (4)
43	woorden waarin /ks/ geschreven wordt als x	–	3 (3)
44	verkleinwoorden -aatje, -eetje, -ootje, -uutje en met de uitgang -nkje	4 (4)	3 (3)
45	woorden met assimilatieverschijnselen	3 (3)	3 (4)

De digitale toets M7 werkwoorden bestaat, net als de papieren toets M7 werkwoorden, voor iets minder dan de helft uit opgaven die ook in toets E7 werkwoorden voorkomen. Ook hier komen de aantallen opgaven per categorie grotendeels overeen met die van de papieren variant. Bij de toets M7 werkwoorden zien we verschillen bij de categorieën 2.a, 2.b en 2.c, bij de toets E7 niet-werkwoorden bij de categorieën 1.b, 2.b, 2.d, 3.b, 3.c en 4.a. De afwijking bestaat elke keer uit één opgave meer of minder ten opzichte van de papieren toets. De precieze aantallen opgaven per categorie staan in tabel 3.1b.

Tabel 3.1b Spellingcategorieën in de digitale toetsen Spelling 3.0 groep 7 werkwoorden  
(met de aantallen in papieren toetsen tussen haakjes)

	Categorie	Omschrijving	M7	E7
1. o.t.t.	1.a	-t achter stam van zwak ww dat in o.v.t. de uitgang -de krijgt	7 (7)	4 (4)
	1.b	wel of geen -t achter een stam op -d	6 (6)	5 (6)
2. o.v.t.	2.a	zwak ww dat in o.v.t. de uitgang -te(n) of -de(n) krijgt	7 (6)	4 (4)
	2.b	verdubbeling d of t bij zwak ww met stam op -d of -t	6 (7)	5 (6)
	2.c	geen -t bij sterk ww dat in 2e en 3e persoon eindigt op -d	6 (7)	6 (6)
	2.d	uitgang -sde(n) of -fde(n) bij zwak ww met stam op -z of -v	7 (7)	6 (5)
3. voltooid deelwoord	3.a	keuze voor eind-d of eind-t bij zwakke werkwoorden met een stam die <b>niet</b> eindigt op -d, -t, -v of -z	–	5 (5)
	3.b	homofone gevallen	–	5 (6)
	3.c	zwakke werkwoorden met stam op -d, -t, -v of -z	–	5 (4)
4. (on)voltooid deelwoord bijvoeglijk gebruikt	4.a	wel of geen -n aan het eind; -d of -t aan het eind; onvoltooid deelwoord bijvoeglijk gebruikt	–	5 (4)

Net als de papieren toetsen bevatten de digitale toetsen opgaven van uiteenlopende moeilijkheidsgraad. De toetsen zijn hierdoor geschikt om verschillen tussen leerlingen in beeld te brengen. Een goede illustratie hiervan en van de samenstelling van de digitale toetsen zijn de figuren in bijlage 1: p50- en p80-kanspunten van de opgaven in de digitale toetsen voor groep 7 in relatie tot de gemiddelde vaardigheidsscore voor de afnamemomenten. In deze figuren is de verdeling van de opgaven over de taken van de toetsen visueel weergegeven. De balkjes in de figuren geven het p50- (onderkant van het balkje) en p80-kanspunt (bovenkant van het balkje) van elke opgave aan. Het p50-punt geeft de vaardigheidsscore aan waarbij er sprake is van een kans van 50% om een opgave goed te beantwoorden. In deze figuren is zichtbaar dat de toetsen opgaven bevatten van uiteenlopende moeilijkheidsgraad.

Bij de digitale toetsen Spelling niet-werkwoorden M7 en E7 zijn er makkelijke opgaven (die liggen onder de stippellijn van M7 cq. E7), opgaven van gemiddelde moeilijkheid (doorkruisen de lijn van M7 cq. E7) en enkele moeilijke opgaven (liggen boven de lijn van M7 cq. E7) opgenomen. De meeste opgaven hebben een gemiddelde moeilijkheid. Ook zijn er naar verhouding veel opgaven relatief gemakkelijk, zodat de leerlingen een prettige toetservaring beleven.

Een vergelijkbaar beeld is te zien bij de digitale toetsen Spelling werkwoorden M7 en E7. Ook deze toetsen hebben een spreiding van makkelijke opgaven, opgaven van gemiddelde moeilijkheid en moeilijke opgaven. Wel is het aandeel relatief moeilijke opgaven iets groter dan bij de digitale toetsen Spelling niet-werkwoorden. Ditzelfde beeld zagen we bij de papieren toetsen.

In de tabellen 3.2a en 3.2b zijn de ranges en de gemiddelden weergegeven voor de p-waarden en de  $r_{it}$ -waarden van de items van de digitale toetsen M7 en E7. Bij alle toetsen is te zien dat de p-waarden liggen tussen 0,39 en 0,90. Er is gestreefd naar p-waarden van de items tussen 0,40 en 0,90. Twee items (in E7 werkwoorden) hebben een p-waarde net onder de 0,40. Er is gestreefd naar een goede spreiding van moeilijkheid over de items. De gemiddelde moeilijkheid van de M7 en E7 voor niet-werkwoorden en werkwoorden ligt tussen 0,65 en 0,69.

Bij vier opgaven ligt de  $r_{it}$ -waarde onder 0,20 (M7 werkwoorden). De gemiddelde  $r_{it}$ -waarde is voor alle vier de digitale toetsen 0,34 of hoger. Door de Cotan wordt een  $r_{it}$ -waarde hoger dan 0,30 gekwalificeerd als goed. Met een gemiddelde van 0,34 of hoger is de itemkwaliteit van de toetsen goed te noemen.

Bijlage 2 bevat een volledig overzicht van de p-waarden en de  $r_{it}$ -waarden van de items van de digitale toetsen.

Voor de volledigheid hebben we ook de ranges en de gemiddelden weergegeven voor de p-waarden en de  $r_{it}$ -waarden van de items van de papieren toetsen M7 en E7 voor zowel niet-werkwoorden als werkwoorden (zie de tabellen 3.3a en 3.3b die overgenomen zijn uit de wetenschappelijke verantwoording van de papieren toetsen). Te zien is dat de digitale toetsen iets moeilijker zijn dan de papieren toetsen. Dit heeft als consequentie dat een leerling in een digitale toets iets minder goed hoeft te hebben dan in een papieren toets van hetzelfde niveau om eenzelfde vaardigheidsscore te behalen.

*Tabel 3.2a Range en gemiddelde van p- en  $R_{it}$ -waarden voor de digitale toetsen M7 en E7 van Spelling 3.0 niet-werkwoorden*

Toets	p-waarden		Rit-waarden		N items
	Range	Gemiddelde	Range	Gemiddelde	
M7 nww	0,40 – 0,89	0,69	0,22 – 0,58	0,42	50
E7 nww	0,48 – 0,90	0,69	0,26 – 0,59	0,44	50

*Tabel 3.2b Range en gemiddelde van p- en  $R_{it}$ -waarden voor de digitale toetsen M7 en E7 van Spelling 3.0 werkwoorden*

Toets	p-waarden		Rit-waarden		N items
	Range	Gemiddelde	Range	Gemiddelde	
M7 ww	0,40 – 0,86	0,65	0,17 – 0,60	0,35	40
E7 ww	0,39 – 0,88	0,67	0,20 – 0,57	0,34	50

*Tabel 3.3a Range en gemiddelde van p- en  $R_{it}$ -waarden voor de papieren toetsen M7 en E7 van Spelling 3.0 niet-werkwoorden*

Toets	p-waarden		Rit-waarden		N items
	Range	Gemiddelde	Range	Gemiddelde	
M7	0,35 – 0,90	0,71	0,26 – 0,60	0,47	50
E7	0,41 – 0,88	0,70	0,31 – 0,59	0,46	50

*Tabel 3.3b Range en gemiddelde van p- en  $R_{it}$ -waarden voor de papieren toetsen M7 en E7 van Spelling 3.0 werkwoorden*

Toets	p-waarden		Rit-waarden		N items
	Range	Gemiddelde	Range	Gemiddelde	
M7	0,43 – 0,92	0,70	0,25 – 0,59	0,38	40
E7	0,41 – 0,91	0,71	0,24 – 0,61	0,39	50

Hoewel de digitale toetsen dus niet identiek zijn aan de papieren toetsen, maakt het voor de resultaten niet uit of leerlingen de papieren of de digitale toetsen maken. De papieren en de digitale opgaven konden namelijk door middel van papier-digitaal vergelijkingsonderzoek in één opgavenbank niet-werkwoorden dan wel werkwoorden, ondergebracht worden; dat wil zeggen dat ze één en dezelfde vaardigheidsschaal (niet-werkwoorden cq. werkwoorden) representeren. Elke verzameling opgaven, of dit nu digitale opgaven



zijn of opgaven op papier, is dan geschikt om die vaardigheid te toetsen, mits de betreffende verzameling min of meer is afgestemd op het niveau van de doelgroep. Zie voor een verantwoording hoofdstuk 4.

### 3.3 Statistische beschrijving

Voor het schatten van de vaardigheid van een leerling maken we bij de LVS-toetsen in het algemeen gebruik van twee berekeningswijzen. Bij de eerste berekeningswijze wordt uitgegaan van het aantal goed op de toets en worden de opgaven niet gewogen met de discriminatie-index: elke opgave telt even zwaar mee in de berekening. Deze berekeningswijze maakt gebruik van de zogenoemde ongewogen score. Bij de tweede berekeningswijze worden de opgaven wél gewogen met hun discriminatie-index, er is dan sprake van gewogen scores. Statistisch gezien is de tweede berekeningswijze te prefereren: de gewogen score (bij gebruik van OPLM) is namelijk een voldoende statistiek voor de (latente) vaardigheid. Met andere woorden, alle informatie over de vaardigheid kunnen we bepalen met behulp van de gewogen score. Voor de schattingswijze van de (latente) vaardigheid waarbij gebruik gemaakt wordt van de ongewogen score geldt dat we (een klein beetje) informatie verliezen. Bovendien geldt dat de schattingen asymptotisch identiek zijn én dat beide schattingen gelijk zijn aan de ware vaardigheid. Het gebruik van de gewogen score heeft één voordeel boven het gebruik van de ongewogen score: de standaardfout van de schatting is aanzienlijk kleiner. Een nadeel is echter dat voor het berekenen van de gewogen score het gehele antwoordpatroon van de leerling nodig is. Dat vraagt voor scholen die de papieren toetsen handmatig verwerken een grote tijdsinspanning. Bij de digitale toetsen speelt dit nadeel niet, omdat de antwoorden op de computer automatisch worden opgeslagen en verwerkt. Daarom wordt bij het schatten van de vaardigheidsscores van digitale toetsen altijd gebruikgemaakt van de gewogen scores.

Voor de verschillende boekjes (per afnamemoment) in het kalibratieonderzoek papier-digitaal zijn de correlaties tussen de schattingen met de beide genoemde berekeningswijzen allen groter dan 0,99. Er zijn voor deze situatie T-testen uitgevoerd: voor alle leerlingen geldt dat er geen significant verschil is tussen beide schattingen. Voor de beide reeksen (per afnamemoment) is ook gekeken naar het teken van de verschillen: vaardigheidsscore gewogen papier versus vaardigheidsscore gewogen digitaal én vaardigheidsscore gewogen digitaal versus vaardigheidsscore digitaal ongewogen. In beide gevallen is het teken ongeveer even vaak positief als negatief. Voor de vergelijking van de gewogen versus de ongewogen vaardigheidsscores geldt bijvoorbeeld dat in ongeveer de helft van de gevallen de schatting op basis van de gewogen score kleiner is dan die op basis van de ongewogen scores. We kunnen hieruit concluderen dat het voor het schatten van de vaardigheid geen verschil maakt welke van de beide schattingsmethoden (gewogen of ongewogen) wordt gebruikt. Omdat het bij digitale toetsen makkelijk te implementeren is om de gewogen score te gebruiken en omdat deze theoretisch preciezer is, wordt bij digitale toetsen de voorkeur gegeven aan de gewogen score.

In de tabellen 3.4a en 3.4b is te zien dat de vaardigheidsverdelingen exact hetzelfde zijn als die van de verdelingen van de papieren versie. De reden hiervoor is dat voor de normering van de digitale toetsen de normen van de papieren toetsen zijn aangehouden. De gegevens zijn gebaseerd op 2342 leerlingen voor M7 niet-werkwoorden, 2579 leerlingen voor E7 niet-werkwoorden, 540 leerlingen voor M7 werkwoorden en 2418 leerlingen voor E7 werkwoorden. Dit betreft de aantallen leerlingen van de normering van de papieren toetsen. De waarden laten zien dat de vaardigheidsverdeling bij benadering normaal is. De figuren 3.1 tot en met 3.4 met de verdeling van de vaardigheidsscores laten dit ook zien. Ook deze figuren zijn precies zo opgenomen in de wetenschappelijke verantwoording van de papieren toetsen Spelling 3.0 groep 7. Voor de duidelijkheid geven we deze figuren op de volgende pagina's nogmaals weer.

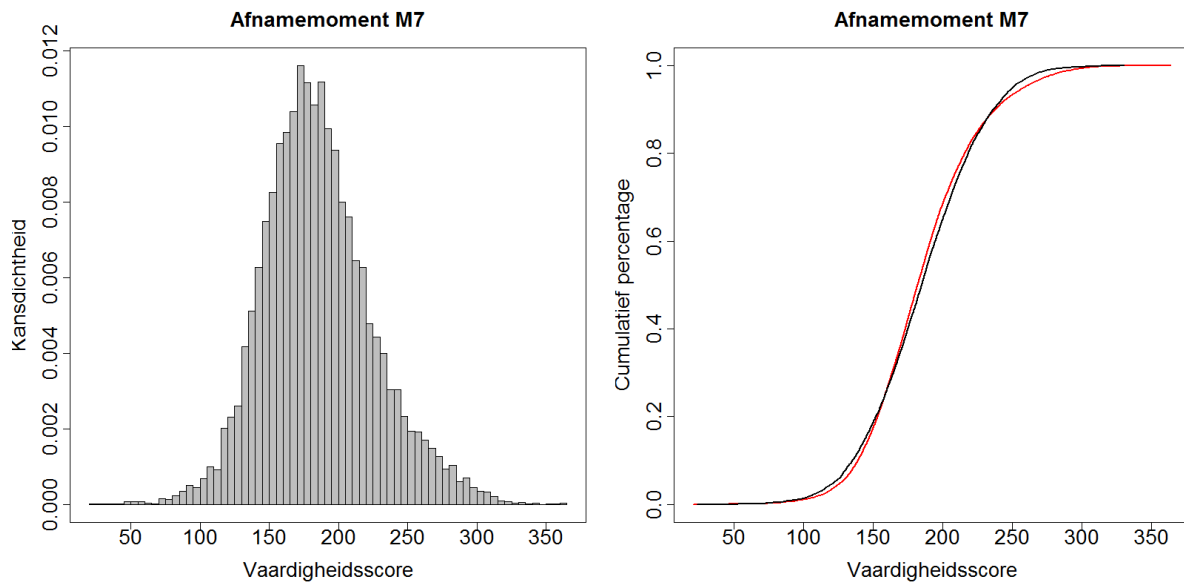
Tabel 3.4a Beschrijvende gegevens digitale toetsen Spelling niet-werkwoorden M7 en E7 op de gewogen scoreschaal en de vaardigheidsschaal

	Gemiddelde	Standaarddeviatie	Kurtosis	Scheefheid
M7 nww gewogen score	133,5	40,6	-0,37	-0,63
M7 nww vaardigheid	351,3	32,3	0,63	0,41
E7 nww gewogen score	138,2	43,6	-0,32	-0,70
E7 nww vaardigheid	356,7	28,5	0,29	-0,06

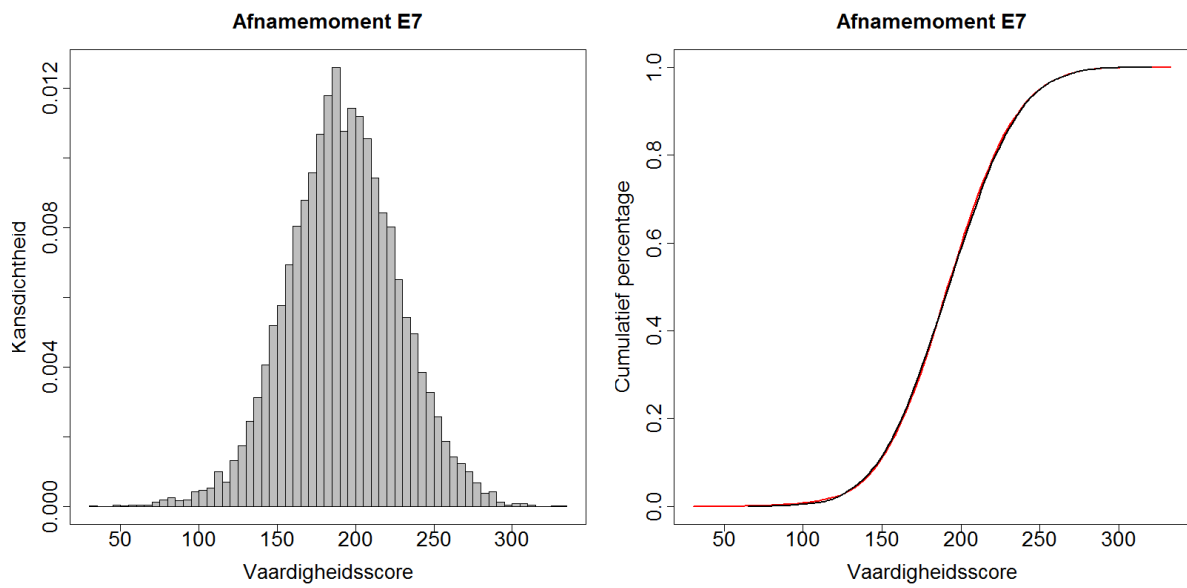
Tabel 3.4b Beschrijvende gegevens digitale toetsen Spelling werkwoorden M7 en E7 op de gewogen scoreschaal en op de vaardigheidsschaal

	Gemiddelde	Standaarddeviatie	Kurtosis	Scheefheid
M7 ww gewogen score	95,7	34,9	-0,98	-0,17
M7 ww vaardigheid	131,7	23,7	0,16	0,10
E7 ww gewogen score	115,7	78,5	-1,62	-0,26
E7 ww vaardigheid	136,2	23,2	0,65	0,25

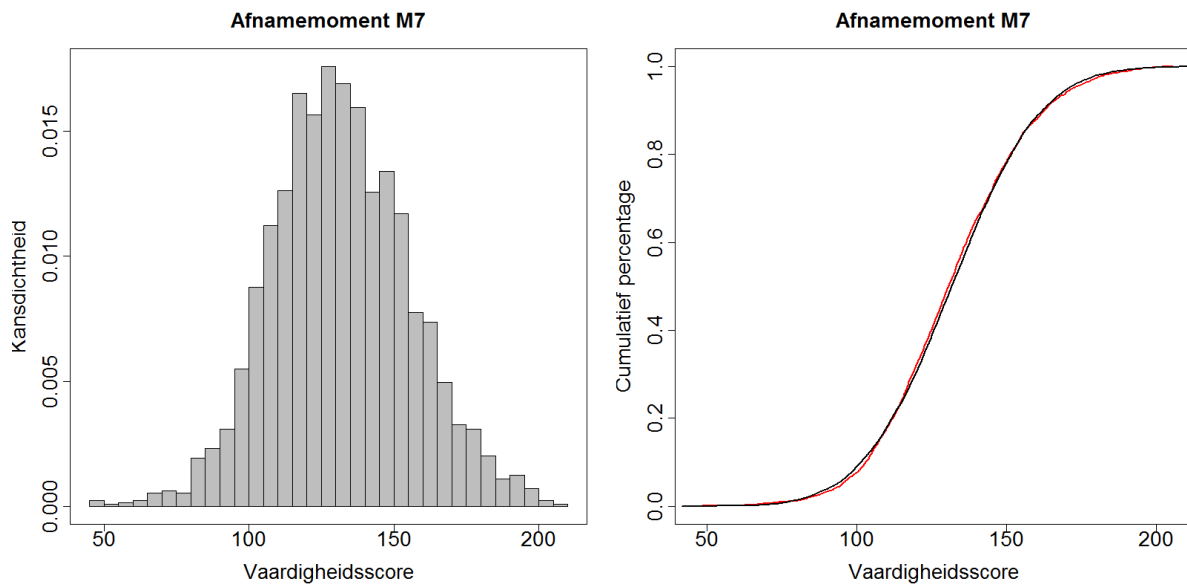
Figuur 3.1 Weergave van behaalde vaardigheidsscores en cumulatieve verdeling M7 niet-werkwoorden



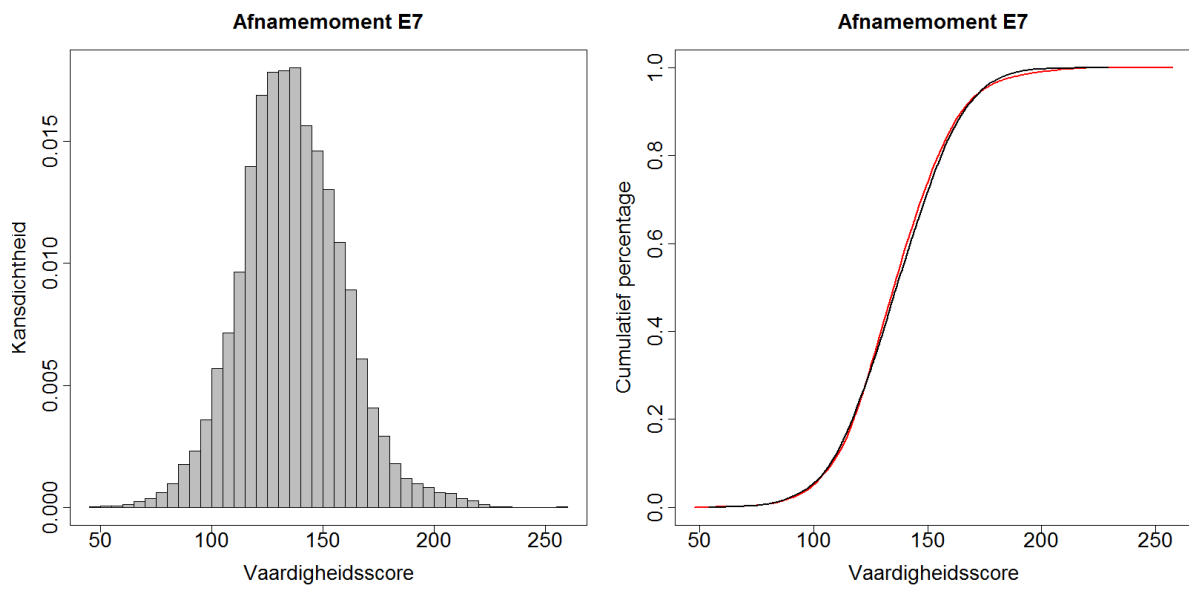
Figuur 3.2 Weergave van behaalde vaardigheidsscores en cumulatieve verdeling E7 niet-werkwoorden



Figuur 3.3 Weergave van behaalde vaardigheidsscores en cumulatieve verdeling M7 werkwoorden



Figuur 3.4 Weergave van behaalde vaardigheidsscores en cumulatieve verdeling E7 werkwoorden



## 4 Kalibratie en normering

### 4.1 Rationale van de kalibratieonderzoeken

Aan het begin van het toetsontwikkelingsproces van de LVS-toetsen Spelling 3.0 groep 7, zijn in 2012 en 2013 opgaven geconstrueerd. Deze opgaven zijn allereerst in *papieren* versie onderzocht in kalibratieonderzoeken in 2014 en 2015 en – na opname van de opgaven in de uit te geven ‘papieren’ toetsen – in normeringsonderzoeken in 2016 en 2017. Zie voor het kalibratie- en normeringsonderzoek van de papieren items voor groep 7 paragraaf 4.2 van de Wetenschappelijke verantwoording van de LVS- toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018).

In 2017 en 2018 hebben vervolgens kalibratieonderzoeken voor de *digitale* items plaatsgevonden. Dit gebeurde in de vorm van papier-digitaal vergelijkingsonderzoek. De hoofdvraag in deze kalibratieonderzoeken was: meten de papieren items en de digitale items dezelfde vaardigheid? We kunnen deze vraag bevestigend beantwoorden als de papieren en digitale items op dezelfde schaal blijken te passen. Hieronder in paragraaf 4.2 beschrijven we deze papier-digitaal vergelijkingsonderzoeken en rapporteren we de resultaten. De conclusie is dat de digitale items dezelfde vaardigheid meten als de papieren items. Dit betekent dat de papieren en de digitale items op één schaal passen, dat de papieren en digitale versies van toetsen van hetzelfde niveau (bijvoorbeeld M7) leiden tot dezelfde vaardigheidsschatting en dat dezelfde normering gehanteerd kan worden. Dit geldt zowel voor Spelling niet-werkwoorden als voor Spelling werkwoorden. De normering lag al vast voor de papieren toetsen. We nemen hieronder in paragraaf 4.3 de normeringsgegevens in verkorte versie over van de verantwoording van de papieren toetsen, omdat er sprake is van één normering die zowel geldt voor de papieren toetsen als de digitale toetsen.

In dit hoofdstuk zal worden besproken hoe de digitale en papieren opgaven op de vaardigheidsschaal spelling niet-werkwoorden dan wel werkwoorden passen. De papieren en de digitale opgaven vormen daarmee ook één opgavenbank. Dat betekent onder andere dat de vaardigheid van een leerling met elke willekeurige selectie van opgaven, ook een selectie van *digitale* opgaven, uit deze bank gemeten kan worden. Hoe nauwkeurig dat gebeurt hangt uiteraard af van het aantal opgaven en van de psychometrische eigenschappen van de opgaven, onder andere de moeilijkheid van de gekozen opgaven in relatie tot de vaardigheid van de leerling. Daarmee is meteen ook het tweede doel van de hier gerapporteerde analyses gegeven: het onderbouwen van de keuze van de items die – gegeven de doelgroep – het beste in de digitale toetsen passen. We kunnen laten zien dat de digitale toetsen voor groep 7 uiteindelijk opgaven bevatten met psychometrisch goede eigenschappen (zie de tabellen 5.1a en 5.1b). Deze psychometrische eigenschappen komen overeen met die van de papieren versies. We kunnen met de digitale opgaven dus even goed de vaardigheid spelling meten als met de papieren opgaven. Er is daarom geen onderzoek nodig naar de equivalentie van de scores op de verschillende versies: de scores zijn immers per definitie uitwisselbaar. Dit betekent ook dat de normen voor de papieren en de digitale toetsen hetzelfde kunnen en moeten zijn. Immers, met beide toetsen meten we dezelfde vaardigheid bij eenzelfde populatie.

### 4.2 Kalibratieonderzoek digitale items

#### 4.2.1 Opzet van het kalibratieonderzoek voor de digitale items

In aparte kalibratieonderzoeken is onderzocht of de digitale items bij de papieren items op de schaal Spelling niet-werkwoorden dan wel werkwoorden passen. In januari 2017 vond het papier-digitaal-kalibratieonderzoek M7 niet-werkwoorden plaats. In juni 2017 vond het papier-digitaal-kalibratieonderzoek E7 niet-werkwoorden plaats, tegelijkertijd met het papier-digitaal-kalibratieonderzoek E7 werkwoorden. In januari 2018 ten slotte vond het papier-digitaal-kalibratieonderzoek M7 werkwoorden plaats.

Alle opgaven van de papieren toetsen 3.0 groep 7 zijn omgezet naar digitale versies, aangevuld met reserve-opgaven. In een papier-digitaal onderzoek is het design onvolledig: in tegenstelling tot een volledig papieren onderzoek overlappen de boekjes slechts gedeeltelijk. De echte link om de boekjes op één schaal te krijgen, verloopt via de papieren uitgave. In de tabellen met afnamedesigns (tabel 4.1 t/m 4.4) zijn de boekjes van de papieren uitgave ook opgenomen. Hier is te zien dat de overlap via de papieren uitgave ervoor zorgt dat de afnamedesigns verbonden zijn. Alle nieuwe (digitale) taken zijn immers afgenomen bij leerlingen die ook de papieren Starttaak van LVS Spelling tweede generatie gemaakt hebben. Omdat de spellingvaardigheid van een leerling niet verandert tijdens een toets, kan de vaardigheid op het papieren gedeelte van de toets vergelijkbaar worden geacht met de vaardigheid op het digitale gedeelte van de toets. Hierdoor zijn de twee afnamemethoden verbonden. Voor de digitale items zijn wel eigen moeilijkheids- en discriminatieparameters geschat. Want ook al zijn de items van de papieren starttaak en de digitale starttaak inhoudelijk gelijk, door de verschillende afnamemethoden kunnen ze niet beschouwd worden als dezelfde items. Het is immers niet noodzakelijkerwijs zo dat de papieren versie en digitale versie van een item precies even moeilijk zijn en/of even goed discrimineren. We bespreken hieronder opzet en design voor de kalibratieonderzoeken.

#### *Medio 7 niet-werkwoorden*

In het papier-digitaal onderzoek voor het 'medio' (M) afnamemoment van januari 2017 zijn 100 items voorgelegd aan 429 leerlingen van groep 7. Elke leerling maakte één digitale taak met 35 nieuwe items uit LVS Spelling 3.0: ofwel taak 1 3.0 ofwel taak 2 3.0. De taken 1 en 2 bestonden elk uit 25 gedigitaliseerde items van de papieren toets M7 3.0, aangevuld met 10 reserve-items. Daarnaast maakte elke leerling de Starttaak M7 van LVS tweede generatie, bestaande uit 30 items. Ongeveer de helft van de leerlingen maakte de papieren Starttaak van de tweede generatie (de leerlingen die toegewezen waren aan boekje 1 en 2) en ongeveer de helft maakte de digitale Starttaak van de tweede generatie (de leerlingen die toegewezen waren aan boekje 3 en 4). Er is voor gezorgd dat geen enkele leerling beide versies (papier en digitaal) van eenzelfde item kreeg voorgelegd.

Doordat de helft van de onderzochte groep de Starttaak van tweede generatie op papier maakte, konden we de papieren en digitale items vergelijken: passen ze op een en dezelfde schaal? Doordat ook een grote groep leerlingen beide taken digitaal maakte, waren we in staat de beste *digitale* items te selecteren voor de definitieve digitale toetsen.

De items waren verdeeld over vier verschillende opgavenboekjes in een onvolledig, maar 'verbonden' design. De data van het papier-digitaal onderzoek werden gekoppeld aan de data die verzameld werden in het normeringsonderzoek voor de papieren uitgave M7 uit 2016. Om de koppeling te realiseren werden de data die verzameld zijn in het papier-digitaal onderzoek toegevoegd aan de dataset van het genoemde normeringsonderzoek (die dient voor de schaling van de items in de itembank Spelling).

Tabel 4.1 Afnamedesign proefonderzoek papier-digitaal M7 niet-werkwoorden

Boekje	Taak				Aantal leerlingen
	Papier	Digitaal			
	Starttaak M7 LVS 2 <sup>e</sup> generatie	Starttaak M7 LVS 2 <sup>e</sup> generatie	Taak 1 3.0	Taak 2 3.0	
1					127
2					95
3					88
4					119
Aantal leerlingen per taak	222	207	215	214	

### Eind 7 niet-werkwoorden

In het papier-digitaal onderzoek voor het 'einde' (E) afnamemoment van juni 2017 zijn 100 items voorgelegd aan 669 leerlingen van groep 7. Ook hier maakte elke leerling één digitale taak met 35 nieuwe items uit LVS Spelling 3.0: ofwel taak 1 3.0 ofwel taak 2 3.0. De taken 1 en 2 bestonden elk uit 25 gedigitaliseerde items van de papieren toets E7 3.0, aangevuld met 10 reserve-items. Daarnaast maakte elke leerling de Starttaak E7 van LVS tweede generatie, bestaande uit 30 items. Ongeveer de helft van de leerlingen maakte de papieren Starttaak van de tweede generatie (de leerlingen die toegewezen waren aan boekje 1 en 2) en ongeveer de helft maakte de digitale Starttaak van de tweede generatie (de leerlingen die toegewezen waren aan boekje 3 en 4). Er is voor gezorgd dat geen enkele leerling beide versies (papier en digitaal) van eenzelfde item kreeg voorgelegd.

Doordat de helft van de onderzochte groep de Starttaak van de tweede generatie op papier maakte, konden we de papieren en digitale items vergelijken: passen ze op een en dezelfde schaal? Doordat ook een grote groep leerlingen beide taken digitaal maakte, waren we in staat de beste *digitale* items te selecteren voor de definitieve digitale toetsen.

De items waren verdeeld over vier verschillende opgavenboekjes in een onvolledig, maar 'verbonden' design. De data van het papier-digitaal onderzoek werden gekoppeld aan de data die verzameld werden in het normeringsonderzoek voor de papieren uitgave E7 uit 2016. Om de koppeling te realiseren werden de data die verzameld zijn in het papier-digitaal onderzoek toegevoegd aan de dataset van het genoemde normeringsonderzoek (die dient voor de schaling van de items in de itembank Spelling).

Tabel 4.2 Afnamedesign proefonderzoek papier-digitaal E7 niet-werkwoorden

Boekje	Taak				Aantal leerlingen
	Papier	Digitaal			
	Starttaak E7 LVS 2 <sup>e</sup> generatie	Starttaak E7 LVS 2 <sup>e</sup> generatie	Taak 1 3.0	Taak 2 3.0	
1					169
2					175
3					144
4					181
Aantal leerlingen per taak	344	325	313	356	

### Medio 7 werkwoorden

De opzet van het papier-digitaal onderzoek voor het 'medio' (M) afnamemoment groep 7 wijkt af van de andere papier-digitaal onderzoeken. Dit komt omdat er geen toets Spelling werkwoorden M7 van de tweede generatie was. Verder valt op dat het papier-digitaal onderzoek van M7 werkwoorden later plaatsvindt dan die van E7. De reden hiervoor is dat pas later besloten is om ook een toets Spelling werkwoorden te maken voor M7, waardoor het onderzoek - net als het papieren normeringsonderzoek - later plaatsvond dan het onderzoek voor E7.

In het papier-digitaal onderzoek M7 van januari 2018 zijn 50 items voorgelegd aan 350 leerlingen van groep 7. Elke leerling maakte twee taken met nieuwe items uit LVS Spelling 3.0. Twee groepen leerlingen maakten zowel een taak op papier als een taak digitaal. Een derde groep leerlingen maakte beiden taken digitaal. De taken 1 en 2 bestonden elk uit 20 items van papieren toets M7 werkwoorden, aangevuld met 5 reserve-items. De papieren taken 1 en 2 bevatten dezelfde items als de digitale taken 1 en 2.

Doordat de meeste leerlingen zowel papieren items als digitale items maakten, konden we de papieren en digitale items vergelijken: passen ze op een en dezelfde schaal? Doordat ook een groep leerlingen beide taken digitaal maakte, waren we in staat de beste *digitale* items te selecteren voor de definitieve digitale toets. De items waren verdeeld over drie verschillende opgavenboekjes in een onvolledig, maar ‘verbonden’ design. De data van het papier-digitaal onderzoek werden gekoppeld aan de data die verzameld werden in het normeringsonderzoek voor de papieren uitgave M7 werkwoorden uit 2017. Om de koppeling te realiseren werden de data die verzameld zijn in het papier-digitaal onderzoek toegevoegd aan de dataset van het genoemde normeringsonderzoek (die dient voor de schaling van de items in de itembank Spelling).

Tabel 4.3 Afnamedesign proefonderzoek papier-digitaal M7 werkwoorden

Boekje	Taak				Aantal leerlingen
	Papier		Digitaal		
	Taak 1 3.0	Taak 2 3.0	Taak 1 3.0	Taak 2 3.0	
1					152
2					76
3					122
Aantal leerlingen per taak	152	76	198	274	

#### Eind 7 werkwoorden

In het papier-digitaal onderzoek voor het ‘einde’ (E) afnamemoment van juni 2017 zijn 85 items voorgelegd aan 613 leerlingen van groep 7. Elke leerling maakte één digitale taak met 30 nieuwe items uit LVS Spelling 3.0: ofwel taak 1 3.0 ofwel taak 2 3.0. De taken 1 en 2 bestonden elk uit 25 gedigitaliseerde items van de papieren toets E7 3.0, aangevuld met 5 reserve-items. Daarnaast maakte elke leerling de Starttaak E7 werkwoorden van LVS tweede generatie, bestaande uit 25 items. Ongeveer de helft van de leerlingen maakte de papieren Starttaak van de tweede generatie (de leerlingen die toegewezen waren aan boekje 1 en 2) en ongeveer de helft maakte de digitale Starttaak van de tweede generatie (de leerlingen die toegewezen waren aan boekje 3 en 4). Er is voor gezorgd dat geen enkele leerling beide versies (papier en digitaal) van eenzelfde item kreeg voorgelegd.

Doordat de helft van de onderzoeksgroep de Starttaak van de tweede generatie op papier maakte, konden we de papieren en digitale items vergelijken: passen ze op een en dezelfde schaal? Doordat ook een grote groep leerlingen beide taken digitaal maakte, waren we in staat de beste *digitale* items te selecteren voor de definitieve digitale toetsen.

De items waren verdeeld over vier verschillende opgavenboekjes in een onvolledig, maar ‘verbonden’ design. De data van het papier-digitaal onderzoek werden gekoppeld aan de data die verzameld werden in het normeringsonderzoek voor de papieren uitgave E7 werkwoorden uit 2016. Om de koppeling te realiseren werden de data die verzameld zijn in het papier-digitaal onderzoek toegevoegd aan de dataset van het genoemde normeringsonderzoek (die dient voor de schaling van de items in de itembank Spelling).



Tabel 4.4 Afnamedesign proefonderzoek papier-digitaal E7 werkwoorden

Boekje	Taak				Aantal leerlingen
	Papier	Digitaal			
	Starttaak E7 LVS 2 <sup>e</sup> generatie	Starttaak E7 LVS 2 <sup>e</sup> generatie	Taak 1 3.0	Taak 2 3.0	
1					177
2					213
3					118
4					105
Aantal leerlingen per taak	390	223	295	318	

#### 4.2.2 De stappen in de kalibratie

Met kalibratie wordt bedoeld dat we kengetallen zoeken bij de items die de antwoorden van de leerlingen goed representeren. Hoe de kengetallen gezocht worden, ligt deels vast door het gekozen model. Hoe succesvol deze operatie is, kan statistisch getoetst worden. Eenvoudig gezegd, schatten we in OPLM met de CML-methode de itemparameters en controleren we of deze de data goed voorspellen. Voor een exacte beschrijving van de statistische toetsen die in OPLM gebruikt worden, hun eigenschappen en feitelijke implementatie in OPLM, verwijzen we naar Verhelst (1993). Hier beperken we ons tot een korte beschrijving van de principes van de statistische toetsen die gebruikt zijn in de kalibratieprocedure. De statistische toetsen in OPLM hebben goede statistische en asymptotische eigenschappen, daar OPLM behoort tot de exponentiële familie, met de gewogen somscore:

$$S = \sum_{i=1}^k a_i x_i, \quad (4.1)$$

als een 'afdoende statistiek' (*sufficient statistic*) voor de vaardigheid  $\theta$ . Dit betekent dat alle informatie in de data met betrekking tot de vaardigheid in deze statistiek aanwezig is. Hiervan wordt gebruikgemaakt bij de statistische toetsen in OPLM. Het basisprincipe van de statistische toetsen in OPLM is dat op grond van de afdoende statistiek  $S$  de personen in de data kunnen worden gegroepeerd. En binnen deze groepen kan de verwachte proportie goede antwoorden op een item onder het model,  $p(+|s)$ , vergeleken worden met de feitelijk geobserveerde proportie goede antwoorden,  $prop(+|s)$ . Via de basisvergelijking van OPLM kunnen we eenvoudig de conditionele kans op het goed beantwoorden van de items afleiden en daarmee kunnen we  $p(+|s)$  evalueren,  $prop(+|s)$  volgt uit de data. Discrepancies tussen  $p(+|s)$  en  $prop(+|s)$  duiden op schendingen van het model. Deze discrepanties vormen de basis voor de diverse statistische toetsen in OPLM. De toetsingsgrootheid voor de veronderstelde discriminatie-indices is gegeven door

$$M = f_{s \in H} (p(+|s) - prop(+|s)) + f_{s \in L} (prop(+|s) - p(+|s)). \quad (4.2)$$

Deze zogenoemde M-toetsen verdelen de scoregroepen in een laag deel ( $L$ ) en een hoog deel ( $H$ ) en  $f$  is een monotone functie. M-toetsen hebben een duidelijke interpretatie: is M significant positief dan is de veronderstelde steilheid van de ICC (item karakteristieke curve) overschat in het model, is M daarentegen erg laag dan is de index te klein. Verhelst laat zien voor welke functie,  $f$ ,  $M \approx N(0,1)$ . In OPLM zijn drie verschillende M-toetsen geïmplementeerd die verschillen in de definitie van de hoge en lage scoregroepen. Naast deze M-toetsen is er een algemene itemtoets die de volgende vorm heeft

$$S = f(p(+|s) - prop(+|s)).$$

Deze zogeheten S-toets heeft een  $\chi^2$  verdeling onder het model. Als globale modeltoets is de R1c-toets

(Glas, 1988) geschikt. Ook de distributie van alle afzonderlijke S-toetsen komt hiervoor in aanmerking. Als we deze S-toetsen opvatten als onafhankelijk, wat ze strikt genomen niet zijn, dan zouden de overschrijdings-kansen uniform verdeeld moeten zijn op het (0,1) interval.

Kortom, als we afzien van de formeel-statistische achtergrond van de gehanteerde toetsen, kan de kalibratieprocedure als volgt worden samengevat:

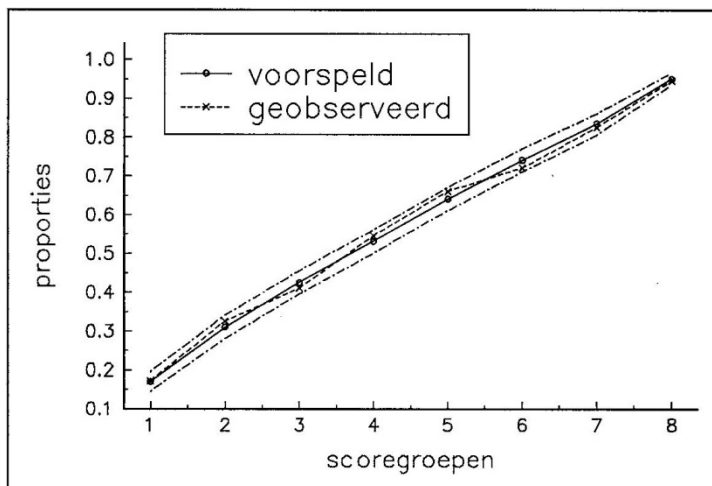
1. Met behulp van het programma OPCAT stellen we de discriminatie-indices in OPLM in en hercoderen we indien noodzakelijk de antwoordcategorieën in de data.
2. Vervolgens schatten we de itemparameters met behulp van de CML-methode.
3. Met behulp van de M-toetsen controleren we of de discriminatie-indices goed zijn ingesteld.
4. Een volgende controle betreft de overschrijdingskansen van de S-toetsen en een grafische modelcontrole door middel van het programma WOPPLOT (grafische inspectie van de ICC's).
5. Tot slot vindt een globale modelcontrole plaats in de vorm van een R1c-toets en de verdeling van de overschrijdingskansen van de S-toetsen.

De stappen 1 tot en met 5 worden een aantal malen doorlopen tot het resultaat bevredigend is. Afhankelijk van de uitkomsten kunnen items worden verwijderd. Ook inhoudelijke overwegingen spelen een rol in dit beslissingsproces.

#### 4.2.3 Toetsing van het IRT-model

Het is niet eenvoudig om de kwaliteit van de kalibratie aan te tonen. De belangrijkste statistische instrumenten om de passing van een opgave in het IRT-model te bewerkstelligen en uiteindelijk te documenteren betreffen de hierboven al besproken S-toetsen. Het lastige daarvan is, dat de toetsing voor een groot deel visueel gebeurt. Dit kunnen we illustreren aan de hand van figuur 4.1 (zie Staphorsius, 1994, blz. 239). Figuur 4.1 beeldt voor een opgave de gegevens af waarop de betreffende S-toetsen gebaseerd zijn (zie handleiding OPLM: Verhelst, 1992). Ten behoeve van deze toetsing wordt de totale groep van leerlingen die een verzameling opgaven gemaakt heeft, ingedeeld in een aantal (meestal acht) scoregroepen. Elke groep bestaat uit leerlingen met een ongeveer even hoge score. De geobserveerde proporties juiste antwoorden van deze groepen (telkens gesymboliseerd door een x) zijn door de middelste stippellijn verbonden. De volle lijn daarentegen verbindt de proporties die op grond van de parameterschattingen voorspeld kunnen worden. De twee buitenste lijnen geven het 95%-betrouwbaarheidsinterval aan. De breedte van dit interval is in belangrijke mate afhankelijk van het aantal leerlingen dat de opgave heeft beantwoord. Uit de figuur blijkt heel duidelijk dat de geobserveerde proporties, zoals bedoeld, binnen het 95%-betrouwbaarheidsinterval van de (geschatte) voorspelde proporties liggen, en dit komt in grote lijnen overeen met een niet-significante S-toetsingsgrootte (Verhelst et al., 1994).

Figuur 4.1 Grafische voorstelling van een  $S_i$ -toets



Het is ondoenlijk om voor alle opgaven dergelijke grafische voorstellingen in deze verantwoording op te nemen. Daarom beperken we ons steeds per digitale toetsversie tot het item met de slechtste en de beste S-passing, aangevuld met een qua S-toetsingsresultaat gemiddelde (dat wil zeggen, meest representatieve) passing. De voorbeelden in de figuren 4.2a en 4.2b illustreren dat voor de toetsen voor groep 7 zelfs bij de slechtst passende opgave sprake is van een zeer aanvaardbaar beeld. Er wordt in dit geval voor een deel (van de onderscheiden scoregroepen) niet beantwoord aan de eis dat de geobserveerde proportie binnen het 95%-betrouwbaarheidsinterval van de geschatte proporties ligt. Dit beeld doet zich slechts bij enkele opgaven voor die dan ook een uitzondering vormen. De overige opgaven voldoen voor alle scoregroepen wel aan die eis. De afbeeldingen voor de representatieve en best passende opgaven illustreren dit. Dit leidt tot de conclusie dat bij vrijwel alle opgaven in de toetsen Spelling een grafische voorstelling van de S-toetsing hoort die in grote lijnen met figuur 4.1 overeenkomt. Dit is een zeer sterke aanwijzing dat het meetinstrument en het meetmodel dat ontwikkeld is, respectievelijk gebruikt is, adequaat zijn om het gedrag van de leerlingen te verklaren. Bovendien blijkt, en dat is vanuit theoretisch oogpunt nog belangrijker, dat gemeten verschillen in gedrag tussen de leerlingen te verklaren zijn door één unidimensionaal concept. Dat laatste wordt nog beter duidelijk voor de voorbeelden van S-toetsen van de 'totale' kalibraties van groep 7. Hierin wordt nog beter duidelijk hoe goed parameters van de items die uitsluitend op papier afgenomen zijn passen bij de parameters die digitaal zijn afgenomen. Er is nauwelijks tot geen sprake van differentieel itemfunctioneren (DIF). Met de 'totale' kalibraties wordt overigens de volledige opgavenbank van papieren en digitale items van groep 7 bedoeld. Meer over deze 'totale' kalibraties staat in paragraaf 4.2.4.

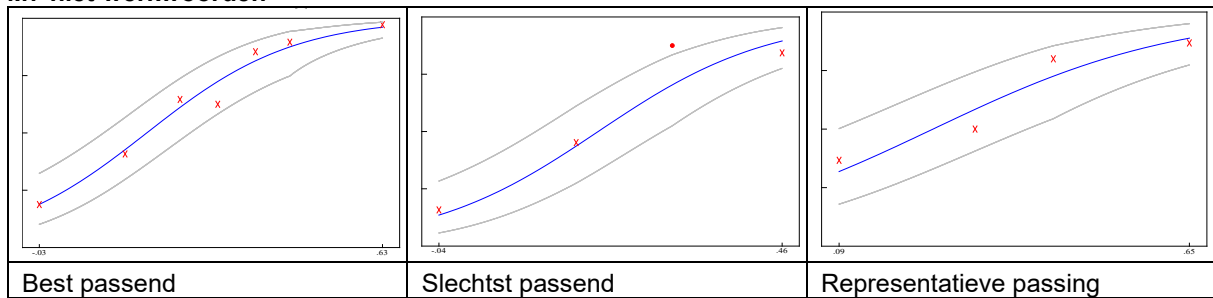
#### 4.2.4 Totale kalibratie per groep

Om goed te kunnen bepalen of de digitale items (en hun parameters) passen bij de papieren parameters en om te bekijken of de items op één schaal passen is ook een 'totale' kalibratie uitgevoerd per groep, voor zowel Spelling niet-werkwoorden als voor Spelling werkwoorden. Dat wil zeggen dat voor alle items van de twee toetsen Spelling niet-werkwoorden van groep 7 een kalibratie werd uitgevoerd, alsook voor alle items van de twee toetsen Spelling werkwoorden van groep 7. Hierbij werden de parameters van de op papier afgenomen items gefixeerd op de waarden zoals geschat in de normeringsonderzoeken van papier. Door vervolgens de parameters van de 'digitale' items op dezelfde kalibratieschaal te schatten als de 'papieren' items kan goed bepaald worden of deze items bij elkaar op een schaal passen. Ook kan bepaald worden of er sprake is van differentieel itemfunctioneren (DIF) van de digitale items ten opzichte van de papieren items. Er is bij een item sprake van DIF als de digitale versie ervan niet op dezelfde schaal kan worden gebracht. Hieronder zal daarom ook de passing van het meetmodel van de totale kalibratie per groep besproken worden.

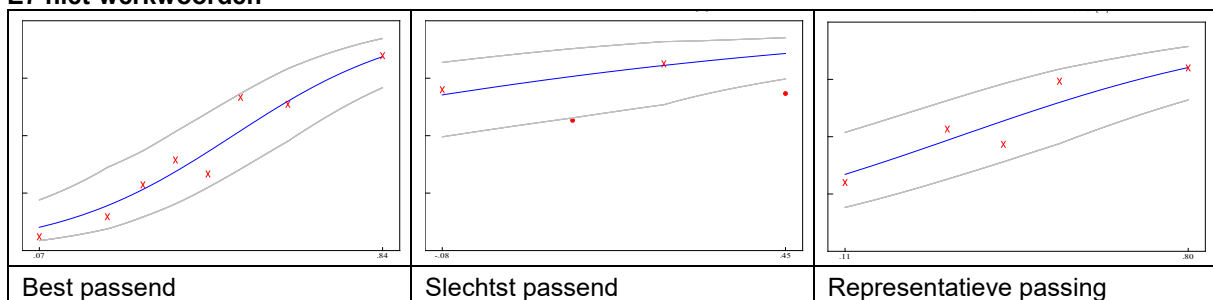
Omdat besloten is om geen aparte normeringen voor de digitale toetsen te ontwikkelen en de normering van de papieren items te gebruiken is de passing van het meetmodel van de totale kalibratie per groep cruciaal om uitspraken te kunnen doen over de kwaliteit van de digitale toetsen LVS Spelling 3.0.

Figuur 4.2a Voorbeelden van S-toetsen voor de digitale toetsen Spelling 3.0 niet-werkwoorden M7 en E7 met per toets de best passende, de slechtst passende en een qua passing representatieve opgave

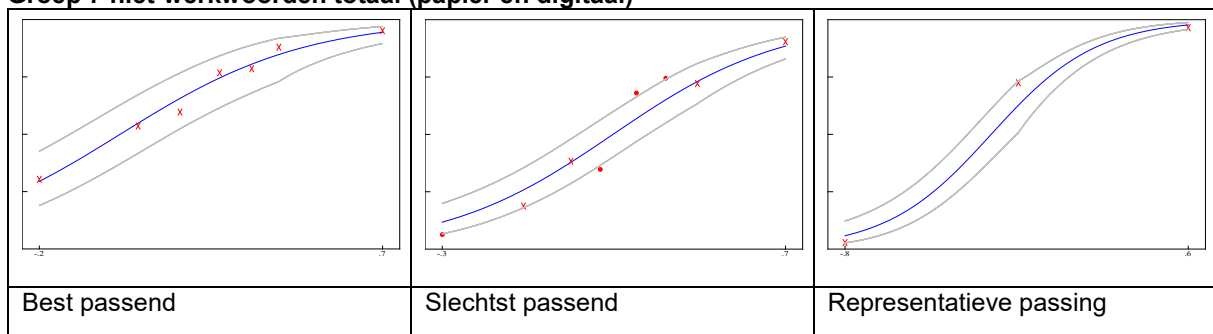
**M7 niet-werkwoorden**



**E7 niet-werkwoorden**

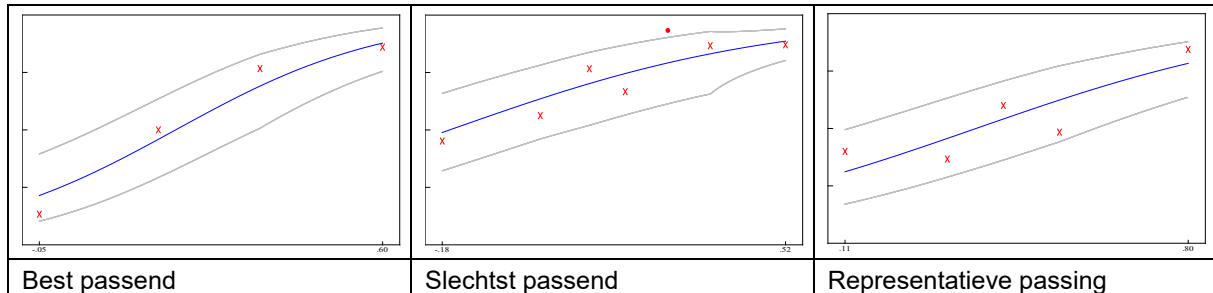


**Groep 7 niet-werkwoorden totaal (papier en digitaal)**

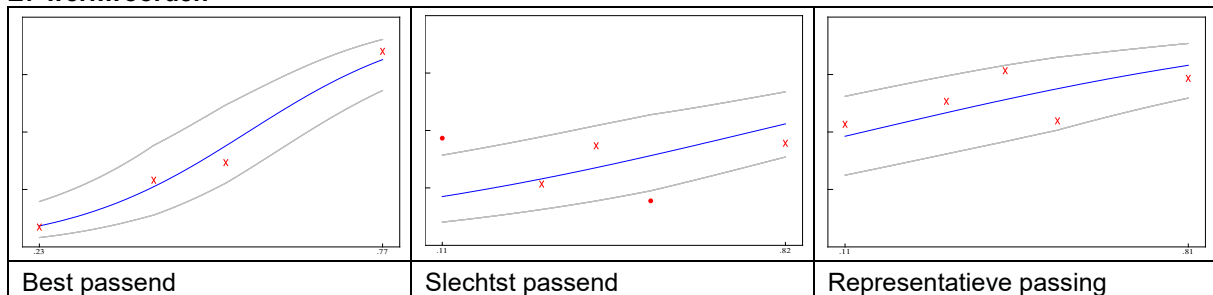


Figuur 4.2b Voorbeelden van S-toetsen voor de digitale toetsen Spelling 3.0 werkwoorden M7 en E7 met per toets de best passende, de slechtst passende en een qua passing representatieve opgave

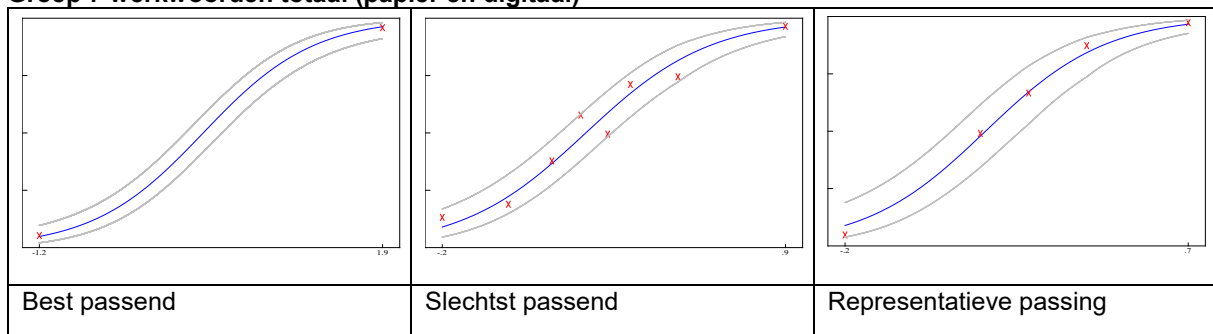
**M7 werkwoorden**



**E7 werkwoorden**



**Groep 7 werkwoorden totaal (papier en digitaal)**



In feite kan men bij de kalibratie beter varen op deze grafische weergaven dan op toetsingsresultaten in termen van exacte getallen en de significantie daarvan. Niettemin zijn er bij de kalibratie S-toetsen uitgevoerd die een indicatie geven van de kwaliteit van de kalibratie. Daarbij zijn we vooral geïnteresseerd in de distributie van de overschrijdingskansen van deze verzameling toetsingsresultaten. Als we de S-toetsen opvatten als onafhankelijk, wat ze strikt genomen niet zijn, dan zouden de overschrijdingskansen uniform verdeeld moeten zijn binnen het (0,1) interval, uiteraard met zo weinig mogelijk significante resultaten. Tabel 4.5a en 4.5b waarin het (0,1) interval is opgedeeld in tien gelijke stukken, geeft een beeld van de uitkomsten bij een kalibratie van alle opgaven van de digitale toetsen Spelling 3.0 groep 7. Daarnaast is aangegeven in hoeveel gevallen de overschrijdingskans kleiner was dan 0,01, respectievelijk 0,05. Het is duidelijk dat voor alle toetsen de verdeling redelijk gelijkmatig is over het gehele interval van overschrijdingskansen. Deze resultaten geven een bevestiging van het eerder geschetste beeld, dat met uitzondering van enkele opgaven, sprake is van niet-significante S-toetsen. Zij vormen een kwantitatieve ondersteuning van de conclusie dat de opgaven een unidimensionaal construct representeren.



Tabel 4.5a Verdeling van overschrijdingskansen bij S-toetsen voor digitale toetsen Spelling 3.0 groep 7 niet-werkwoorden en bij de totale kalibratie

	0	1	2	3	4	5	6	7	8	9	10	
M7 nww	3	5	2	4	3	4	7	3	2	9	4	5
E7 nww	5	2	5	9	6	2	4	7	1	2	2	4
Gr. 7 nww	12	13	29	69	74	63	78	77	71	82	85	115

Tabel 4.5b Verdeling van overschrijdingskansen bij S-toetsen voor digitale toetsen Spelling 3.0 groep 7 werkwoorden en bij de totale kalibratie

	0	1	2	3	4	5	6	7	8	9	10	
M7 ww	1	3	1	2	3	2	3	4	6	5	4	6
E7 ww	1	1	3	4	3	5	4	4	1	3	5	5
Gr. 7 ww	10	19	24	44	37	49	45	49	42	46	35	62

In tabel 4.6a en 4.6b zijn de R1c-waarden weergegeven voor dezelfde afnames waarvoor in tabel 4.5a en 4.5b de resultaten van de S-toetsen zijn weergegeven. R1c is een statistiek die zicht geeft op de modelpassing van de toets als geheel. Voor een goede modelfit geldt als vuistregel dat R1c bij voorkeur niet groter zou moeten zijn dan ongeveer anderhalf maal het aantal vrijheidsgraden (df). Tot tweemaal het aantal vrijheidsgraden geldt het als een acceptabele passing. Uit tabel 4.6a en 4.6b blijkt dat de modelpassing voor de meeste toetsen goed is. Alle R1c-waarden zijn onder of rond anderhalf maal het aantal vrijheidsgraden. De significantie van de statistische toetsingen is bij de grote aantallen in de analyse nauwelijks informatief.

Tabel 4.6a R1c-waarden voor de digitale toetsen Spelling 3.0 niet-werkwoorden M7 en E7

Toetsversie	R1c	df	p
M7 nww	1025,8	849	<0,01
E7 nww	1616,4	1096	<0,01
Gr. 7 nww	13063,8	9098	<0,01

Tabel 4.6b R1c-waarden voor de digitale toetsen Spelling 3.0 werkwoorden M7 en E7

Toetsversie	R1c	df	p
M7 nww	373,1	278	<0,01
E7 nww	1063,9	796	<0,01
Gr. 7 nww	6409,4	4308	<0,01

Ten slotte bespreken we nog een methode om de modelpassing te verantwoorden die wordt besproken in het COTAN Beoordelingssysteem (Evers, Lucassen, Meijer en Sijtsma, 2010, p. 40). Het betreft hier een poging om de nauwkeurigheid van de itemparameterschattingen te beoordelen op basis van een constante (in het COTAN-Beoordelingssysteem met 'c' aangeduid) die weergeeft hoe de relatie is

tussen de standaardfout van de moeilijkheidsparameter van een item en de standaarddeviatie van de vaardigheidsverdeling van de kalibratiepopulatie. Het beoordelingssysteem geeft ook richtlijnen voor het beoordelen van de grootte van deze 'c'. Deze dient te worden beoordeeld als goed als de waarde lager is dan of gelijk aan 0,20. Waarden tussen 0,30 en 0,40 kunnen nog als voldoende worden beschouwd. De waarden voor deze constante zijn weergegeven in tabel 4.7a en 4.7b. De gemiddelde waarden van de constante zijn goed te noemen. Vier items van M7 werkwoorden zijn boven de 0,3 en onder de 0,4. Hetzelfde geldt voor E7 werkwoorden. Van M7 niet-werkwoorden is één item boven de 0,3 en onder de 0,4. Van E7 niet-werkwoorden is geen enkel item boven de 0,3.

De conclusie mag luiden dat we ook op basis van deze analyse de kalibratie geslaagd kunnen noemen.

*Tabel 4.7a Nauwkeurigheid van de itemparameterschattingen (constante 'c') voor digitale toetsen Spelling 3.0 groep 7 niet-werkwoorden*

Toetsversie	Constante 'c'	
	Range	Gemiddelde
M7 nww	0,097 – 0,301	0,160
E7 nww	0,076 – 0,208	0,124

*Tabel 4.7b Nauwkeurigheid van de itemparameterschattingen (constante 'c') voor digitale toetsen Spelling 3.0 groep 7 werkwoorden*

Toetsversie	Constante 'c'	
	Range	Gemiddelde
M7 nww	0,061 – 0,375	0,187
E7 nww	0,065 – 0,364	0,182

Ook op basis van de hierboven beschreven resultaten kan de conclusie luiden dat voor de digitale toetsen Spelling 3.0 M7 en E7, voor zowel niet-werkwoorden als werkwoorden, de kalibratie geslaagd is. Belangrijker nog is de conclusie dat de kalibratie van de schaal waarop de papieren en de digitale items samen gekalibreerd zijn geslaagd is. Hieruit blijkt dat de papieren en de digitale items goed op een schaal passen en er geen sprake is van betekenisvol differentieel itemfunctioneren voor de digitale items in relatie tot de papieren items. Juist deze geslaagde kalibratie maakt het mogelijk om voor de digitale toetsen uit te gaan van de normen van de papieren toetsen.

### 4.3 De normering

De normering die wordt gebruikt voor de digitale toetsen Spelling is gelijk aan de normering van de papieren toetsen Spelling. Dit is mogelijk gezien de koppeling van het papier-digitaal kalibratieonderzoek aan de normeringsonderzoeken voor de papieren uitgaven via het in voorgaande paragraaf besproken design en de geslaagde kalibratie. De (papieren) normering is gebaseerd op de onderliggende (latente) verdeling van de vaardigheid op de afnametijdstippen M7 en E7. Bij de kalibratie is gebleken dat de 'digitale opgaven' op dezelfde schaal geplaatst konden worden als de 'papieren opgaven'. Daardoor kunnen we de eerder gevonden verdelingen van de vaardigheid van de normgroepen (M7 en E7) op deze schaal gebruiken. Voor het beschrijven van de normpopulatie kunnen we daarom gebruikmaken van de eerder gerapporteerde resultaten voor de papieren toetsversies.

De Expertgroep Toetsen PO had als oordeel dat voor de normering van Spelling 3.0 groep 7 een representatieve steekproef is gebruikt. Ook zijn de gebruikte normgroepen groot genoeg en representatief



voor de doelpopulatie, zowel op schoolniveau als leerlingniveau.

Deze normeringen worden in deze paragraaf in verkorte versie besproken. Voor een uitgebreide versie wordt verwezen naar paragraaf 4.3 van de wetenschappelijke verantwoording van de papieren toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018).

Sinds schooljaar 2013/2014 past Cito een nieuwe werkwijze voor het normeren van leerlingvolgsysteemtoetsen toe. Deze werkwijze wordt gebruikt bij het monitoren van de normering van inmiddels uitgegeven toetsen, maar wordt ook gebruikt bij de normering van nieuw uit te geven toetsen, zo ook bij de derde generatie toetsen voor Spelling. De werkwijze die we hieronder beschrijven, komt uit Keuning et al. (2014). Allereerst besteden we aandacht aan de opzet van het normeringsonderzoek, de gehanteerde procedures en de aantallen leerlingen per afnamemoment (paragraaf 4.3.1). Vervolgens komt in paragraaf 4.3.2 de representativiteit van de normsteekproeven aan de orde. De paragraaf wordt afgerond met een presentatie van de resultaten van de normering (i.e. de kenmerken van de vaardigheidsverdelingen op de onderscheiden afnamemomenten; paragraaf 4.3.3).

#### 4.3.1 Opzet

Tijdens het embedded field normeringsonderzoek (zoals omschreven in paragraaf 4.2.1 van de wetenschappelijke verantwoording van de papieren toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018)) werden data verzameld. Om deelnemers te werven voor het normeringsonderzoek zijn scholen aangeschreven. Voor het embedded field normeringsonderzoek is een representatieve steekproef getrokken uit de verzameling van alle basisscholen in Nederland. Dit is gedaan vanuit het bij Cito gebruikelijke steekproefkader dat bepaald wordt door regio, urbanisatiegraad en schooltype (zie verderop voor een omschrijving van deze achtergrondvariabelen).

Voor het normeringsonderzoek M7 niet-werkwoorden zijn in eerste instantie 72 herhalingscholen (scholen die ook meededen aan normeringsonderzoek E6) aangeschreven. Omdat na de eerste aanschrijvingsronde 58% van de herhalingscholen uit normeringsonderzoek E6 bereid bleek deel te nemen aan het normeringsonderzoek, zijn in een tweede aanschrijvingsronde 1166 scholen aangeschreven. Uiteindelijk resulteerde dit in een deelnamebereidheid van zo'n 0,7% van de overige aangeschreven scholen. In totaal meldden zich 50 scholen aan voor het normeringsonderzoek met in totaal 1182 leerlingen. Uiteindelijk was het aantal scholen dat daadwerkelijk gegevens aanleverde gelijk aan 50 en het aantal leerlingen gelijk aan 1180. Daarvan vielen er 10 af vanwege onvolledige gegevens, zodat de antwoorden van 1170 leerlingen konden worden meegenomen in de analyses.

Voor het normeringsonderzoek E7 niet-werkwoorden zijn in totaal 63 herhalingscholen (scholen die ook meededen aan normeringsonderzoek M7) aangeschreven. Omdat na de eerste aanschrijvingsronde 56% van de herhalingscholen uit het normeringsonderzoek M7 ook bereid bleek deel te nemen aan het normeringsonderzoek E7, zijn in een tweede wervingsronde 1977 extra scholen aangeschreven. Uiteindelijk resulteerde dit in een deelnamebereidheid van zo'n 1,4% van de overige aangeschreven scholen (van de tweede werving). In totaal meldden zich 63 scholen aan voor het normeringsonderzoek met in totaal 1375 leerlingen. Uiteindelijk was het aantal scholen dat daadwerkelijk gegevens aanleverde gelijk aan 61 en het aantal leerlingen aan 1311. Daarvan vielen er 25 af vanwege onvolledige gegevens, zodat de antwoorden van 1286 leerlingen konden worden meegenomen in de analyses.

Voor het normeringsonderzoek M7 werkwoorden zijn in totaal 79 herhalingscholen (scholen die ook meededen aan Spelling niet-werkwoorden en/of werkwoorden E7 in juni 2016) aangeschreven en 4161 extra scholen. Na de eerste aanschrijvingsronde bleek 34% van de herhalingscholen ook bereid deel te nemen aan het normeringsonderzoek M7 werkwoorden en 0,1% van de scholen uit de aanvullende steekproef. In totaal meldden zich 33 scholen aan voor het normeringsonderzoek met in totaal 685 leerlingen. Uiteindelijk was het aantal scholen dat daadwerkelijk gegevens aanleverde gelijk aan 28 en het aantal leerlingen 544. Daarvan vielen er 4 af vanwege onvolledige gegevens, zodat de antwoorden van 540 leerlingen konden worden meegenomen in de analyses.

Voor het normeringsonderzoek E7 werkwoorden zijn in eerste instantie 63 herhalingsscholen (scholen die ook meededen aan Spelling niet-werkwoorden M7) aangeschreven. Omdat na de eerste aanschrijvingsronde 60,3% van de herhalingsscholen uit het vorige normeringsonderzoek bereid bleek deel te nemen aan het normeringsonderzoek zijn in een tweede wervingsronde 1977 extra scholen aangeschreven. Uiteindelijk resulteerde dit in een deelnamebereidheid van zo'n 1,2% van de overige aangeschreven scholen. In totaal meldden zich 62 scholen aan voor het normeringsonderzoek met in totaal 1368 leerlingen. Uiteindelijk was het aantal scholen dat daadwerkelijk gegevens aanleverde gelijk aan 60 en het aantal leerlingen gelijk aan 1317. Daarvan vielen er 26 af vanwege onvolledige gegevens, zodat de antwoorden van 1291 leerlingen konden worden meegenomen in de analyses.

Voor het bepalen van de normering werden de gegevens uit het normeringsonderzoek aangevuld met gegevens uit Cito dataretour. In tabel 4.8a en 4.8b zijn de uiteindelijke aantallen scholen en leerlingen in het ('papieren') normeringsonderzoek samengevat.

*Tabel 4.8a Aantal leerlingen per afnamemoment die meegenomen zijn in de normering Spelling niet-werkwoorden*

Afnamemoment	Aantal leerlingen			Aantal scholen
	Steekproef	Dataretour	Normering	Normering
M7	1024*	1318	2342	123
E7	1202*	1377	2579	137

\* Het aantal leerlingen is lager dan dat van de oorspronkelijke normeringssteekproef doordat sommige scholen in de toepassing van het algoritme met gegevens uit Cito dataretour niet zijn geselecteerd.

*Tabel 4.8b Aantal leerlingen per afnamemoment die meegenomen zijn in de normering Spelling werkwoorden*

Afnamemoment	Aantal leerlingen			Aantal scholen
	Steekproef	Dataretour	Normering	Normering
M7	540 (645)	0	540 (645)	27 (32)
E7	1111*	1307	2418	119

\* Het aantal leerlingen is lager dan dat van de oorspronkelijke normeringssteekproef doordat sommige scholen in de toepassing van het algoritme met gegevens uit Cito dataretour niet zijn geselecteerd.

N.B. De leerlingen die in het kalibratieonderzoek papier-digitaal zijn betrokken, maken geen deel uit van de normeringspopulatie.

#### 4.3.2 Representativiteit

Door de werkwijze die werd gevolgd bij de normering is representativiteit van de normeringssteekproeven in principe gegarandeerd. Niettemin werd er een controle uitgevoerd op de representativiteit door de populatieverdelingen verkregen uit gegevens van DUO te vergelijken met de steekproefverdelingen. De steekproef is geanalyseerd in relatie tot de variabelen regio, urbanisatiegraad, schooltype en geslacht. De conclusie is dat de normeringssteekproeven een goede afspiegeling vormen van de populatie. Zie ook Tomesen, Wouda, Krämer & Horsels (2018).

#### 4.3.3 Normeringsresultaten

Na de hierboven beschreven procedure doorlopen te hebben en de normeringssteekproef te hebben samengesteld, kon de normering worden bepaald. Naast het gemiddelde werden de percentielen bepaald. Dat gebeurde op basis van de verdeling van scores die werden gevonden in de normeringssteekproef zoals die is samengesteld op basis van het embedded field normeringsonderzoek en Cito-dataretour. Om de scores van leerlingen te kunnen vergelijken over de tijd worden vaardigheidsscores gebruikt. Uit de ruwe scores van de leerlingen uit het embedded field normeringsonderzoek en Cito-dataretour werden “plausible values” gegenereerd op de nieuw ontwikkelde vaardigheidsschaal. Deze “plausible values” representeren de volledige range aan vaardigheidsscores die een leerling zou kunnen hebben, gegeven de scores. De “plausible values” geven niet alleen informatie over de geschatte vaardigheid maar ook over de onzekerheid die bij die schatting hoort (Keuning et al., 2014). De normering werd vervolgens gebaseerd op de “plausible values” van de leerlingen in de normeringssteekproef. Tabellen 4.9a en 4.9b geven de normgegevens voor de toetsen Spelling 3.0 groep 7.

*Tabel 4.9a Normtabel op leerlingniveau voor Spelling 3.0 groep 7 niet-werkwoorden*

<b>Afname- moment</b>	<b>M</b>	<b>SD</b>	<b>K</b>	<b>S</b>	<b>P10</b>	<b>P20</b>	<b>P25</b>	<b>P40</b>	<b>P50</b>	<b>P60</b>	<b>P75</b>	<b>P80</b>	<b>P90</b>
M7	351,3	32,3	0,63	0,41	314,1	325,5	329,5	341,7	349,0	356,3	370,1	375,8	391,5
E7	356,7	28,5	0,29	-0,06	321,4	333,6	338,4	349,8	356,3	363,6	375,8	380,6	392,8

*Tabel 4.9b Normtabel op leerlingniveau voor Spelling 3.0 groep 7 werkwoorden*

<b>Afname- moment</b>	<b>M</b>	<b>SD</b>	<b>K</b>	<b>S</b>	<b>P10</b>	<b>P20</b>	<b>P25</b>	<b>P40</b>	<b>P50</b>	<b>P60</b>	<b>P75</b>	<b>P80</b>	<b>P90</b>
M7	131,7	23,7	0,16	0,10	103,1	112,3	115,8	124,9	130,7	136,7	147,5	151,6	162,0
E7	136,2	23,2	0,65	0,25	108,2	117,6	121,1	129,6	135,2	140,9	151,0	154,8	164,9

De betreffende normeringstabellen zijn niet alleen van toepassing op de papieren versie van de toetsen Spelling 3.0 voor groep 7, maar ook op de hier verantwoorde digitale versie van deze toetsen.

## 5 Betrouwbaarheid en meetnauwkeurigheid

### 5.1 Betrouwbaarheid

In hoofdstuk 4 is onder meer aangegeven dat elke leerling die deelgenomen heeft aan het normeringsonderzoek slechts een deel van de items gemaakt heeft die uiteindelijk in de toetsen Spelling opgenomen zijn. De betrouwbaarheid van de toetsen in klassieke zin is dan ook niet rechtstreeks te bepalen. Het is echter wel mogelijk om de betrouwbaarheid van iedere toets te schatten door gebruik te maken van het feit dat alle items die zijn opgenomen in de toetsen OPLM-geschaald zijn. Ook andere beschrijvende gegevens, zoals de gemiddelde score en de standaardmeetfout, zijn te schatten op grond van het feit dat de toetsen volledig bestaan uit OPLM-gekalibreerde items. Om relevante beschrijvende gegevens bij de verschillende toetsen te genereren, is gebruikgemaakt van het programma OPLAT (Verhelst, Glas en Verstralen, 1995).

In OPLAT wordt een door Verhelst, Glas en Verstralen (1995, pp. 99-100) ontwikkelde coëfficiënt berekend die qua interpretatie een grote overeenkomst vertoont met betrouwbaarheidscoëfficiënten uit de klassieke testtheorie. Het begrip ware score is wat meer geëxpliciteerd, namelijk als de verwachte score op een (vaste) toets, maar dan gezien als functie van de latente variabele  $\theta$ . Deze verwachte waarde wordt aangeduid met  $\tau(\theta)$ . Als bovendien bekend is hoe  $\theta$  in de populatie verdeeld is, kunnen ook het gemiddelde en de variantie van de ware scores in de populatie bepaald worden. De variantie van de ware scores in de populatie worden aangegeven met het symbool  $Var(\tau)$ . Tussen  $\theta$  en  $\tau(\theta)$  bestaat een een-op-een relatie, immers de een kan uit de andere berekend worden. Het is echter niet zo dat een persoon met vaardigheid  $\theta$  per se de toetsscore  $\tau(\theta)$  moet behalen (dat is alleen zo als de toets oneindig lang wordt).

De geobserveerde score bij een eenmalige afname zal dan ook een afwijking vertonen van de verwachte score, waardoor met een eenmalige toetsafname niet meer zonder fout de waarde van  $\theta$  bepaald kan worden. De variantie van de geobserveerde toetsscore wordt aangegeven met  $Var(t|\tau(\theta))$ , en door weer gebruik te maken van de distributie van  $\theta$  in de populatie kan ook de gemiddelde variantie van de geobserveerde toetsscores berekend worden.

$$Var(t) = E[Var(t | \tau(\theta))] \quad (5.1)$$

Deze variantie kan opgevat worden als de (gemiddelde) meetfoutvariantie in de metriek van de geobserveerde scores ( $t$ ). In analogie met de theorie over de betrouwbaarheid volgt dan

$$MAcc = \frac{Var(\tau)}{Var(\tau) + Var(t)} \quad (5.2)$$

waarin MAcc staat voor Accuracy of Measurement.

Tabellen 5.1a en 5.1b bevatten informatie over de meeteigenschappen van de vaardigheidsschaal Spelling niet-werkwoorden cq. werkwoorden, voor de digitale toetsen voor groep 7. In de eerste kolom staat de aanduiding van de toets. De tweede kolom geeft het aantal items van de toets weer en in de derde kolom staat de maximumscore die gehaald kan worden op de toets. Bij de papieren toetsen is de maximumscore gelijk aan het aantal opgaven dat deel uitmaakt van de totale toets. Bij de digitale toetsen gebruiken we echter de **gewogen** scores, zoals eerder al toegelicht. De vierde kolom geeft de geschatte gemiddelde scores van de leerlingen op de verschillende toetsen. De vijfde kolom bevat informatie over de geschatte standaardmeetfout op de ruwe score van iedere toets. De zesde kolom laat zien wat de geschatte betrouwbaarheidscoëfficiënt (MAcc) van de verschillende toetsen is.

De betrouwbaarheidscoëfficiënten zijn zonder uitzondering hoog. Voor toetsen van het type waar geen zware consequenties voor leerlingen aan verbonden zijn (zoals de toetsen Spelling 3.0 groep 7) geeft de

COTAN aan dat een betrouwbaarheidscoëfficiënt lager dan 0,70 onvoldoende is, een betrouwbaarheidscoëfficiënt tussen 0,70 en 0,80 voldoende, en een betrouwbaarheidscoëfficiënt hoger dan 0,80 goed (Evers, Lucassen, Meijer en Sijsma, 2010). Op grond van dit criterium is de meetnauwkeurigheid van alle toetsen (zeer) goed te noemen.

Tabel 5.1a Beschrijvende gegevens bij de digitale toetsen Spelling 3.0 groep 7 niet-werkwoorden

Toets	Aantal items	Maximum score	Gemiddelde	Standaardmeetfout	MAcc	Test-hertest (simulatie)
M7 nww	50	195	133,5	11,13	0,93	0,93
E7 nww	50	201	138,2	11,35	0,93	0,93

Tabel 5.1b Beschrijvende gegevens bij de digitale toetsen Spelling 3.0 groep 7 werkwoorden

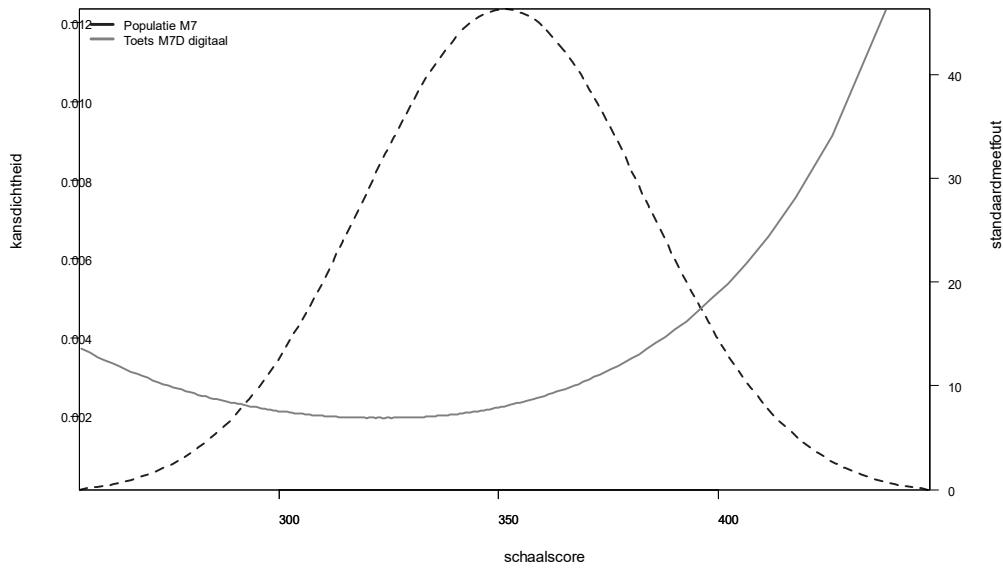
Toets	Aantal items	Maximum score	Gemiddelde	Standaardmeetfout	MAcc	Test-hertest (simulatie)
M7 ww	40	159	95,7	11,50	0,89	0,89
E7 ww	50	202	115,7	8,45	0,98	0,98

Er heeft geen test-hertest onderzoek plaatsgevonden. De afnamecontext van de toetsen Spelling 3.0 leent zich daar niet goed voor. Het feit dat alle items echter OPLM-gekalibreerd zijn, maakt het mogelijk een hertest te simuleren. We hebben een dubbele afname gesimuleerd van 1.000.000 leerlingen. Daarbij hebben we enerzijds de vaardigheidsverdeling van alle leerlingen, anderzijds alle itemparameters als uitgangspunt genomen. Steeds is een bepaalde vaardigheid aselekt uit de verdeling genomen en zijn twee bij deze vaardigheid horende afnames gesimuleerd. Uiteindelijk is de correlatie tussen deze 1.000.000 dubbele (virtuele) afnames berekend. Men kan deze simulatie beschouwen als een test-hertestonderzoek onder ideale condities. De tweede toetsafname is immers volledig onafhankelijk van de eerste toetsafname en wordt niet beïnvloed door de kennis die de leerling mogelijk verworven heeft via de eerste toetsafname. Daarnaast is er geen sprake van invloed van een test-hertest-interval: beide afnames worden gesimuleerd alsof zij op hetzelfde moment plaats zouden vinden. De resultaten zijn weergegeven in de laatste kolom van tabellen 5.1a en 5.1b. De uitkomsten komen vrijwel exact overeen met eerder berekende coëfficiënten en leiden dan ook tot dezelfde conclusies met betrekking tot de betrouwbaarheid van de digitale toetsen Spelling 3.0.

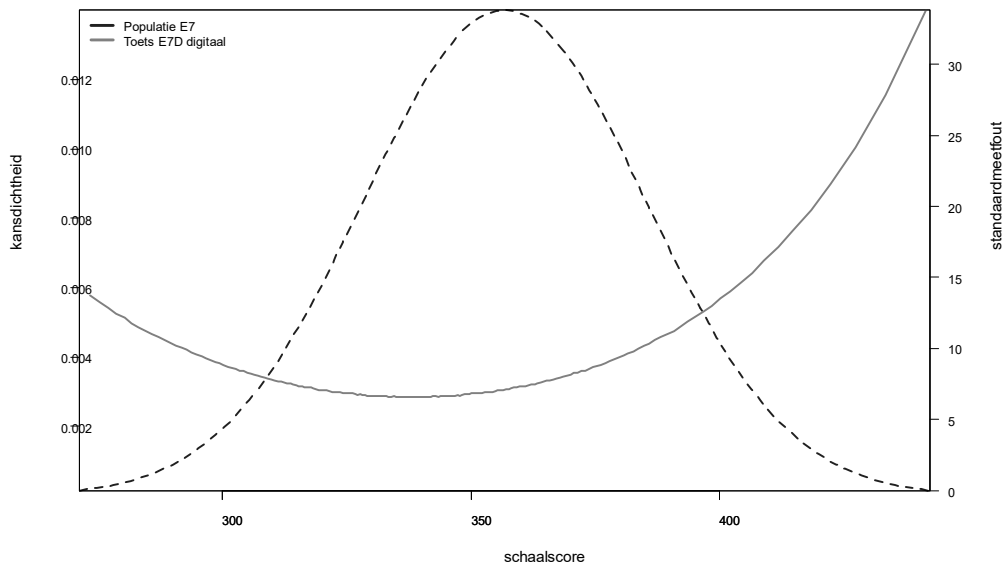
## 5.2 Nauwkeurigheid

De hiervoor vermelde betrouwbaarheidscoëfficiënten hebben alleen betrekking op de globale meetnauwkeurigheid van de toetsen en geven geen beeld van de lokale meetnauwkeurigheid van de digitale toetsen Spelling 3.0. De figuren 5.1 tot en met 5.4 geven grafisch weer hoe het gesteld is met de lokale meetnauwkeurigheid van de verschillende toetsen. In deze figuren staat voor iedere toets de grootte van de meetfout op de vaardigheidsschaal afgebeeld (met verdelingskenmerken zoals aangegeven in tabel 4.9a en 4.9b). Ook zijn de kansdichtheidfuncties voor de normgroepen op de verschillende afnamemomenten opgenomen. Deze laten zien hoe de vaardigheid van de leerlingen verdeeld is over de vaardigheidsschaal in de populatie die de toets gemaakt heeft. De figuren maken duidelijk dat de meetfout kleiner is in de lagere en gemiddelde vaardigheidsregionen dan in de hogere vaardigheidsregionen.

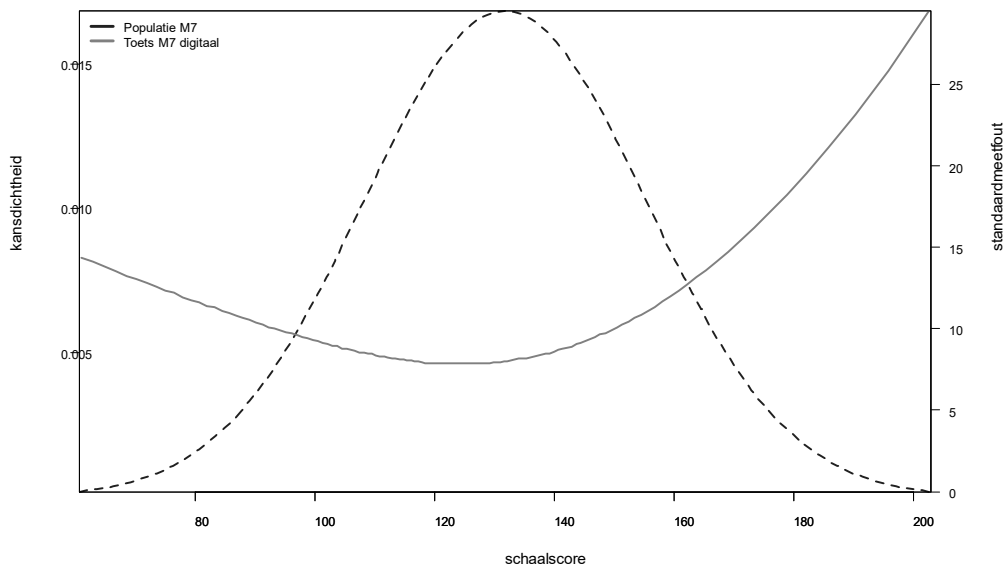
**Figuur 5.1** *Grootte van de meetfouten voor de digitale toets M7 niet-werkwoorden en de kansdichtheidsfunctie voor de M7-populatie*



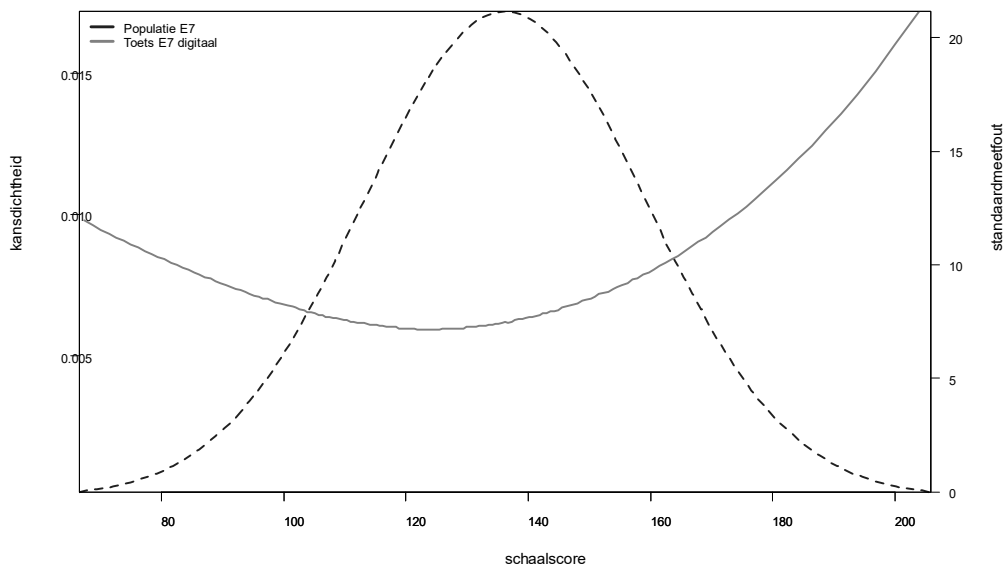
**Figuur 5.2** *Grootte van de meetfouten voor de digitale toets E7 niet-werkwoorden en de kansdichtheidsfunctie voor de E7-populatie*



**Figuur 5.3** Grootte van de meetfouten voor de digitale toets M7 werkwoorden en de kansdichtheidsfunctie voor de M7-populatie



**Figuur 5.4** Grootte van de meetfouten voor de digitale toets E7 werkwoorden en de kansdichtheidsfunctie voor de E7-populatie



### Betrouwbaarheidstabellen

De betekenis van de (lokale) meetnauwkeurigheid voor de beslissingen die met de toets genomen worden, is af te leiden uit betrouwbaarheidstabellen. De tabellen 5.2 tot en met 5.5 laten voor de digitale toetsen M7 en E7, voor zowel niet-werkwoorden als werkwoorden zien hoe vaak de werkelijke vaardigheidsscore in dezelfde scoregroep valt als de geschatte vaardigheidsscore. Zo laat tabel 5.2 zien dat 85,1 procent van de leerlingen die halverwege groep 7 op basis van de digitale M7-toets Spelling niet-werkwoorden in scoregroep V geclassificeerd wordt ook met hun werkelijke vaardigheidsscore in deze scoregroep ingedeeld wordt. De kans dat een V-leerling terecht als V-leerling wordt bestempeld is, met andere woorden, ongeveer 85 procent. Verder laat de linkerkant van tabel 5.2 zien dat 14,7 procent van de leerlingen in scoregroep V een vaardigheidsscore heeft die in werkelijkheid in scoregroep IV valt. De overige getallen in tabellen 5.2 tot en met 5.5 zijn op dezelfde wijze te interpreteren.

Tabel 5.2 Betrouwbaarheidstabel digitale toets M7 niet-werkwoorden voor afnamemoment medio 7

Score- groepen V t/m I	Scoregroep waarin de ware score valt					Score- groepen E t/m A	Scoregroep waarin de ware score valt				
	V	IV	III	II	I		E	D	C	B	A
V	85,1	14,7	0,2	0,0	0,0	E	80,2	19,6	0,2	0,0	0,0
IV	11,9	66,6	20,9	0,6	0,0	D	9,9	67,1	22,9	0,1	0,0
III	0,1	17,0	60,1	22,3	0,6	C	0,0	11,0	70,8	18,1	0,1
II	0,0	0,5	19,6	60,5	19,4	B	0,0	0,0	16,0	67,7	16,2
I	0,0	0,0	0,7	18,5	80,8	A	0,0	0,0	0,2	16,6	83,2

Tabel 5.3 Betrouwbaarheidstabel digitale toets E7 niet-werkwoorden voor afnamemoment einde 7

Score- groepen V t/m I	Scoregroep waarin de ware score valt					Score- groepen E t/m A	Scoregroep waarin de ware score valt				
	V	IV	III	II	I		E	D	C	B	A
V	86,1	13,8	0,1	0,0	0,0	E	81,8	18,1	0,1	0,0	0,0
IV	10,9	68,4	20,2	0,5	0,0	D	9,4	69,4	21,1	0,1	0,0
III	0,1	15,7	61,1	22,6	0,5	C	0,0	10,1	71,6	18,1	0,1
II	0,0	0,4	18,8	61,7	19,0	B	0,0	0,0	14,1	68,3	17,6
I	0,0	0,0	0,6	18,8	80,6	A	0,0	0,0	0,1	15,2	84,6

Tabel 5.4 Betrouwbaarheidstabel digitale toets M7 werkwoorden voor afnamemoment medio 7

Score- groepen V t/m I	Scoregroep waarin de ware score valt					Score- groepen E t/m A	Scoregroep waarin de ware score valt				
	V	IV	III	II	I		E	D	C	B	A
V	80,3	18,9	0,8	0,0	0,0	E	70,6	27,3	2,1	0,0	0,0
IV	16,7	59,2	22,7	1,3	0,0	D	14,0	57,4	28,1	0,5	0,0
III	0,6	21,3	54,5	22,8	0,8	C	0,4	14,8	64,5	20,0	0,3
II	0,0	1,3	21,6	57,4	19,7	B	0,0	0,2	18,2	63,6	18,0
I	0,0	0,0	1,2	19,8	79,1	A	0,0	0,0	0,4	17,4	82,2



Tabel 5.5 Betrouwbaarheidstabel digitale toets E7 werkwoorden voor afnamemoment einde 7

Score-groepen V t/m I	Scoregroep waarin de ware score valt					Score-groepen E t/m A	Scoregroep waarin de ware score valt				
	V	IV	III	II	I		E	D	C	B	A
V	80,6	18,7	0,7	0,0	0,0	E	75,0	24,0	0,9	0,0	0,0
IV	13,8	59,9	24,5	1,8	0,0	D	12,0	60,4	27,2	0,4	0,0
III	0,4	19,7	54,3	24,4	1,1	C	0,2	13,8	65,2	20,5	0,4
II	0,0	1,6	23,0	55,6	19,8	B	0,0	0,2	18,8	62,1	18,9
I	0,0	0,0	1,5	20,7	77,8	A	0,0	0,0	0,4	17,4	82,1

In de onderzoeksliteratuur is weinig geschreven over de beoordeling van betrouwbaarheidstabellen. Wanneer een betrouwbaarheidstabel als goed of voldoende kan worden beschouwd is onduidelijk en wat verwacht mag worden onder ideale omstandigheden is, voor zover ons bekend, niet onderzocht. Daarom worden betrouwbaarheidstabellen vaak samengevat in één of meerdere indices. Wij gebruiken de *plus/minus 1 niveau-index* en de *Marginal Classification Accuracy*. De eerste maat is bedacht door Pilliner (in Wheadon & Stockford, 2011). Hij stelt als ambitieniveau dat 95 procent van de leerlingen in een scoregroep in werkelijkheid ook in die scoregroep moet scoren, **of** één scoregroep daarboven **of** één scoregroep daaronder. In de tabellen zijn dit de gearceerde cellen. Dit ambitieniveau is gebaseerd op de veronderstelling dat geen enkele toets perfect meet en dat er dus altijd sprake is van foutieve classificaties. In dat licht is de maximale accuraatheid die op het individuele niveau bereikt kan worden plus of minus één scoregroep. De tweede maat wordt op verschillende plekken in de literatuur beschreven. De maat laat zien hoe vaak de classificatie op basis van de geschatte vaardigheidsscore gemiddeld gezien overeenstemt met de classificatie op basis van de (gesimuleerde) werkelijke vaardigheidsscore. Bij een ideale, maar in de praktijk niet te realiseren, toetsafname lijkt de *Marginal Classification Accuracy* rond 0,75 - 0,80 uit te komen. In de praktijk liggen de waarden vaak tussen 0,60 en 0,70.

De samenvattende indices voor de toetsen groep 7 zijn te vinden in tabel 5.6a en 5.6b. Waar de betrouwbaarheidstabellen laten zien dat de meeste leerlingen op basis van hun geschatte vaardigheidsscore geplaatst worden in de niveaugroep waar ze werkelijk thuishoren, maken de tabellen 5.6a en 5.6b aannemelijk dat de uitkomsten duidelijk in lijn liggen met het ambitieniveau zoals dat geformuleerd is door Pilliner (in Wheadon & Stockford, 2011) of zelfs boven dit ambitieniveau uitstijgen. Gemiddeld gezien scoort, afhankelijk van de toets en de gekozen indeling in scoregroepen, 98,6 tot 99,9 procent van de leerlingen in een scoregroep ook in werkelijkheid in die scoregroep, **of** één scoregroep daarboven **of** één scoregroep daaronder. De *Marginal Classification Accuracy* loopt uiteen van 65,6 tot 75,1 procent. Dit betekent dat de classificatie op basis van de geschatte vaardigheidsscore gemiddeld gezien in zo'n 70 procent van de gevallen overeenstemt met de classificatie op basis van de (gesimuleerde) werkelijke vaardigheidsscore. De resultaten stemmen hiermee tot grote tevredenheid: het percentage misclassificaties is erg beperkt.

Op basis van bovenstaande gegevens concluderen we dat op basis van de digitale toetsen Spelling 3.0 groep 7 de leerlingen op een betrouwbare manier ingedeeld kunnen worden in normgroepen. Deze indeling voldoet over het algemeen uitstekend gegeven het doel van de toets.

Uiteraard dienen de gebruikers rekening te houden met het gegeven dat er altijd sprake zal zijn van misclassificatie; veelal van maximaal één niveau verschil.

Tabel 5.6a *Samenvattende indices digitale toetsen M7 en E7 niet-werkwoorden*

	Toets M7, afnamemoment M7		Toets E7, afnamemoment E7	
	score- groep I t/m V	score- groep A t/m E	score- groep I t/m V	score- groep A t/m E
Marginal classification accuracy	70,6	73,8	71,6	75,1
Accuracy plus/minus 1 niveau	99,5	99,9	99,5	99,9

Tabel 5.6b *Samenvattende indices digitale toetsen M7 en E7 werkwoorden*

	Toets M7, afnamemoment M7		Toets E7, afnamemoment E7	
	score- groep I t/m V	score- groep A t/m E	score- groep I t/m V	score- groep A t/m E
Marginal classification accuracy	66,1	67,7	65,6	69,0
Accuracy plus/minus 1 niveau	98,8	99,2	98,6	99,5



## 6 Validiteit

Voor de verantwoording van de validiteit verwijzen we naar hoofdstuk 6 van de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018). Alles wat beschreven staat in dit hoofdstuk gaat ook op voor de digitale toetsen Spelling groep 7. Daarbij geldt dat ook de modelfit voor de digitale opgaven bevredigend is (zie paragraaf 4.2.3) en dat daarmee net als bij de papieren versie voldaan wordt aan eisen van unidimensionaliteit als waarborg voor de constructiviteit van de toetsen.



## 7 Samenvatting

In dit hoofdstuk vatten we samen wat in de voorafgaande hoofdstukken is besproken.

In hoofdstuk 1 is aangegeven dat het hier om een *aanvulling* gaat bij de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7 (Tomesen, Wouda, Krämer & Horsels, 2018). Deze aanvulling heeft uitsluitend betrekking op de *digitale* toetsen Spelling 3.0 voor groep 7.

Net als de papieren LVS-toetsen Spelling 3.0 groep 7 vormen de digitale LVS-toetsen Spelling 3.0 voor groep 7 een hulpmiddel om vast te stellen in hoeverre leerlingen kunnen spellen. De toetsen kunnen, in samenhang met de (papieren en digitale) toetsen Spelling 3.0 voor de andere leerjaren, worden gebruikt om de spellingvaardigheid van leerlingen in het primair en speciaal onderwijs in kaart te brengen en om hun ontwikkeling te volgen.

Voor de inhoudelijke aspecten verwijzen we naar de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7. Meetpretentie, gebruiksdoel en functie zijn identiek voor de papieren en digitale toetsen. Specifieke uitgangspunten bij het samenstellen van de digitale toetsen zijn in hoofdstuk 3 beschreven. Dit hoofdstuk bevat ook een beschrijving van enkele psychometrische kenmerken.

In hoofdstuk 4 verantwoorden we de kalibratie- en normeringsonderzoeken. De kalibratieonderzoeken zijn uitgevoerd in de vorm van papier-digitaal vergelijkingsonderzoeken. Hieruit blijkt dat de digitale items op dezelfde schaal passen en daardoor dezelfde vaardigheid meten als de papieren items. De modelfit voor de digitale items is bevredigend en daarmee is voldaan aan de eisen van unidimensionaliteit.

De normering lag al vast voor de papieren items; die hebben we ook aangehouden voor de digitale items.

In hoofdstuk 5 is over de betrouwbaarheidscoëfficiënten gerapporteerd. Net als bij de papieren toetsen, zijn de betrouwbaarheidscoëfficiënten (MAcc's en testhertest) voor de digitale versie van de toetsen voor groep 7 zeer hoog. Ze variëren van 0,89 tot 0,98. Verder zijn in dit hoofdstuk betrouwbaarheidstabellen opgenomen die de betekenis van de meetnauwkeurigheid voor de beslissingen die met de toetsen genomen worden, laten zien. Daarnaast geven we inzicht in de lokale betrouwbaarheid: de meetfout blijkt het kleinst te zijn in de lagere en gemiddelde vaardigheidsregionen.

Voor de verantwoording van de validiteit (hoofdstuk 6) is weer verwezen naar de Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 groep 7.



## 8 Aanvullende literatuur

Cito (2017). *Cito Volgsysteem primair en speciaal onderwijs. Spelling 3.0 Groep 7*. Arnhem: Cito.

Cito (2019). *Spelling 3.0 Handleiding digitale toetsen*. Arnhem: Cito.

Tomesen, M., Wouda, J., Krämer, I. & Horsels, L. (2018). *Wetenschappelijke verantwoording van de LVS-toetsen Spelling 3.0 voor groep 7*. Arnhem: Cito.

Wheadon, C. & Stockford, I. (2011). *Classification accuracy and consistency in GCSE and a level examinations offered by the assessment and qualifications alliance (AQA) november 2008 to june 2009*. Belfast: Office of Qualifications and Examinations Regulation.

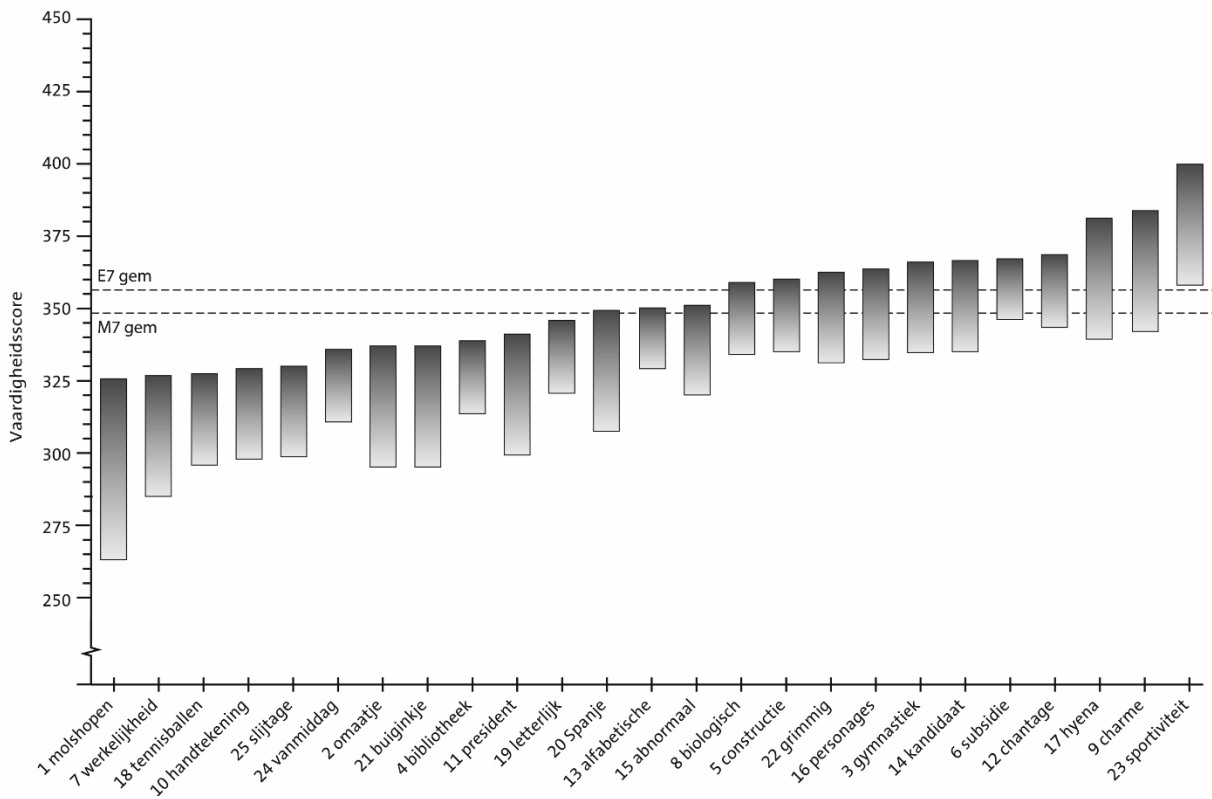




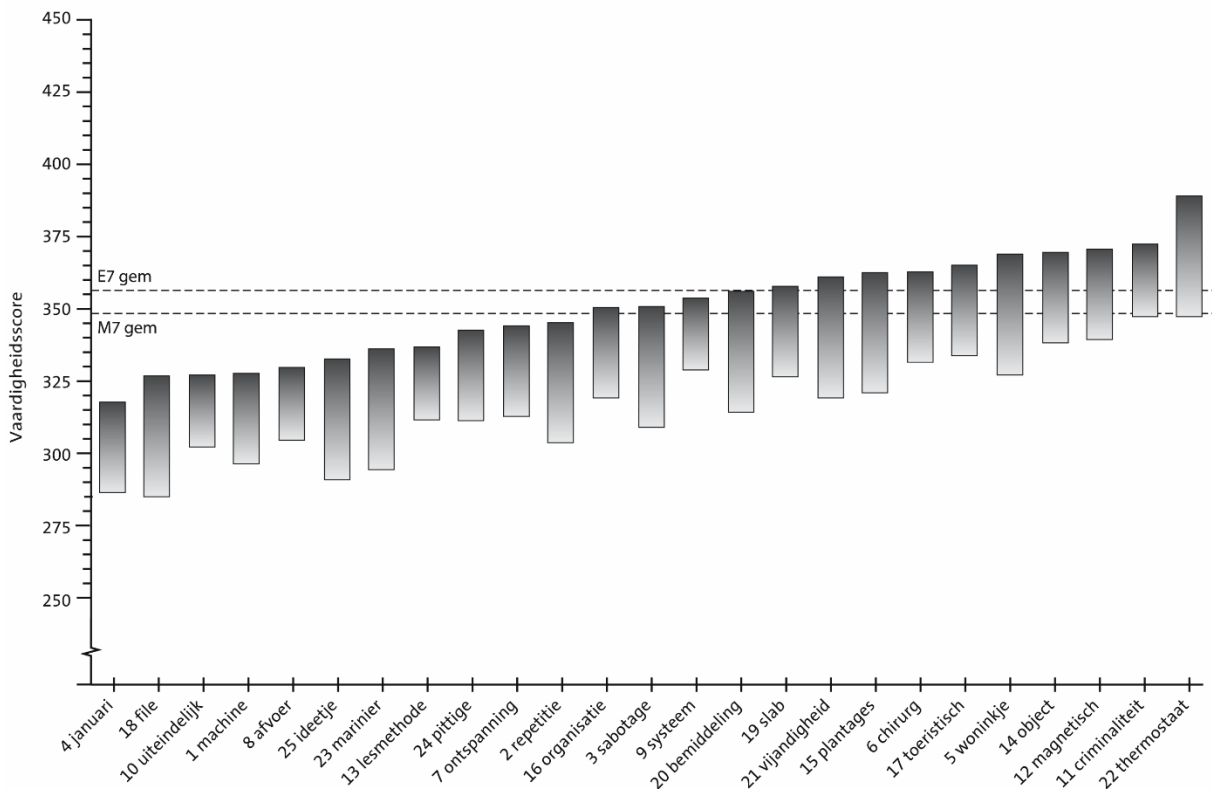
## **Bijlagen**

**Bijlage 1 Moelijkheid van opgaven per taak in Spelling 3.0 digitaal groep 7**

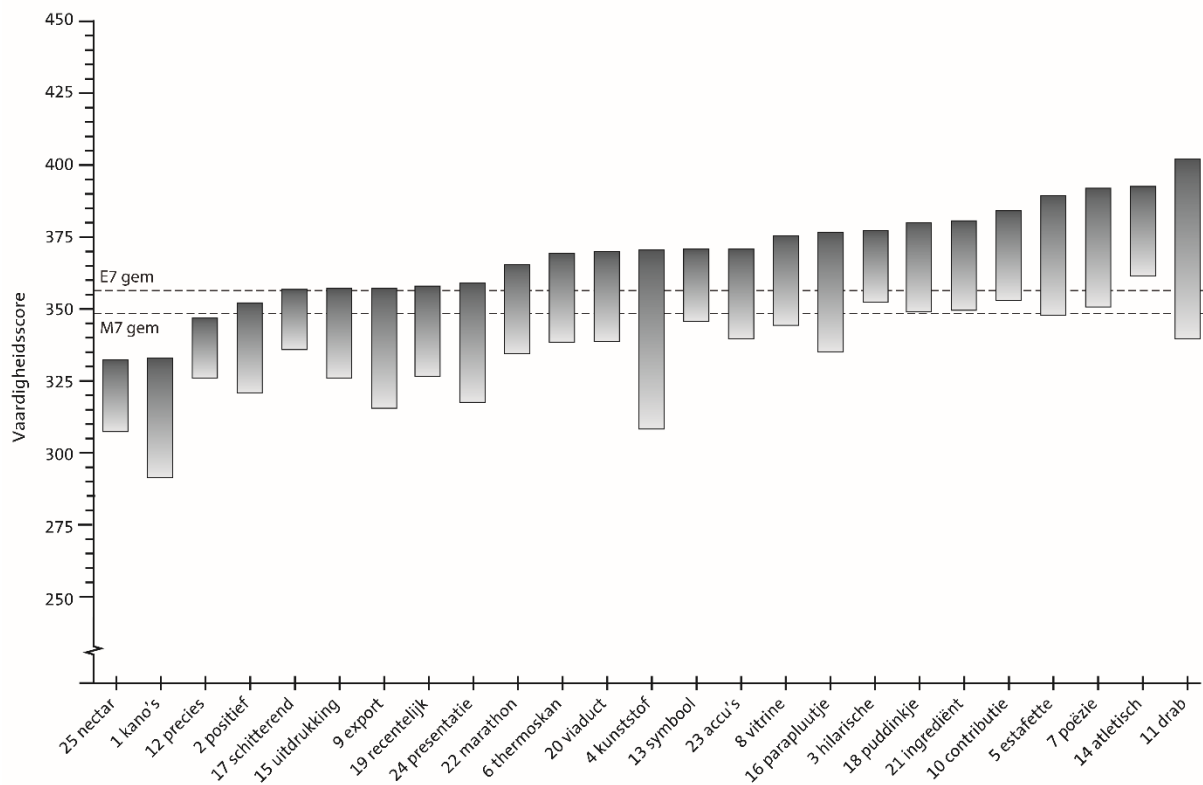
**Spelling niet-werkwoorden M7 digitaal - Taak 1**



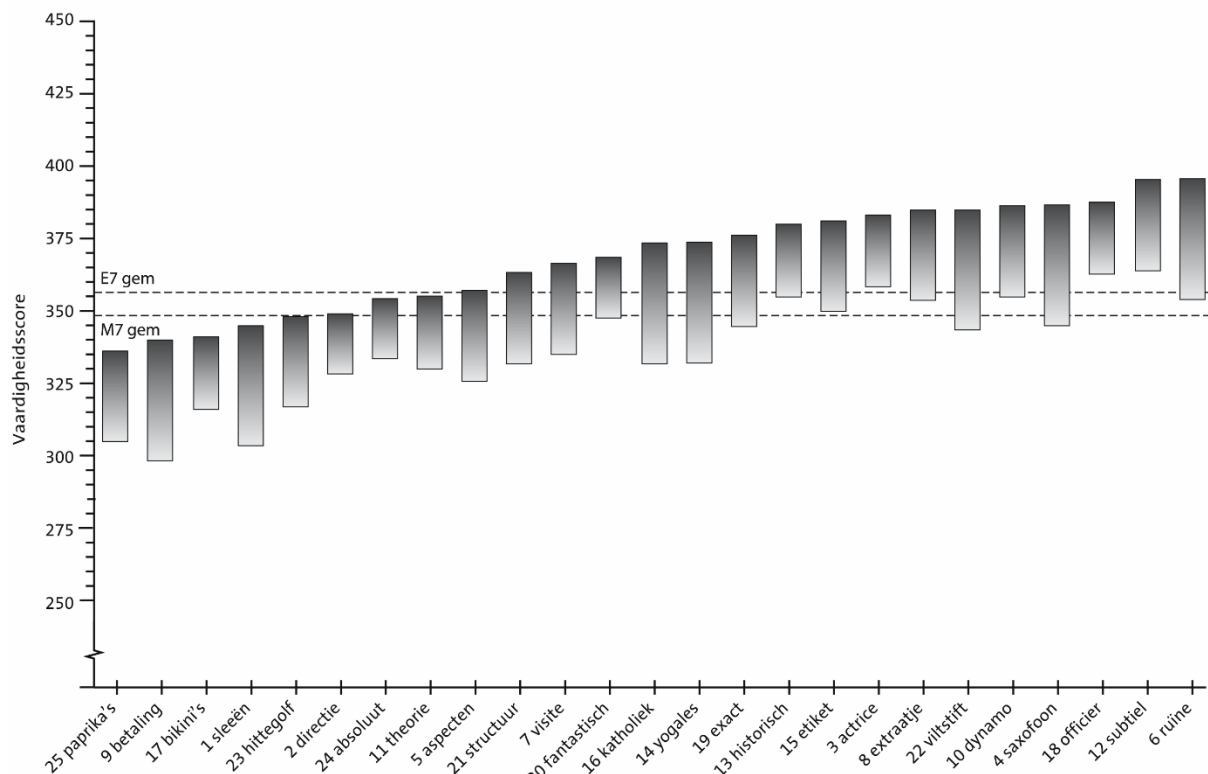
**Spelling niet-werkwoorden M7 digitaal - Taak 2**



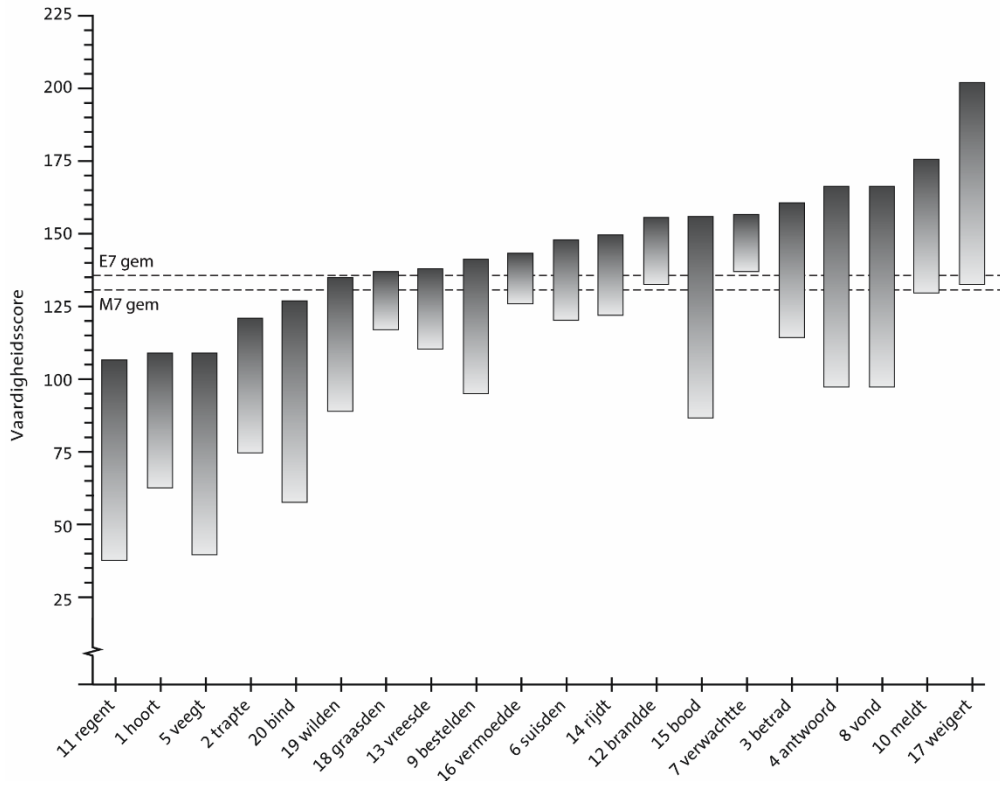
### Spelling niet-werkwoorden E7 digitaal - Taak 1



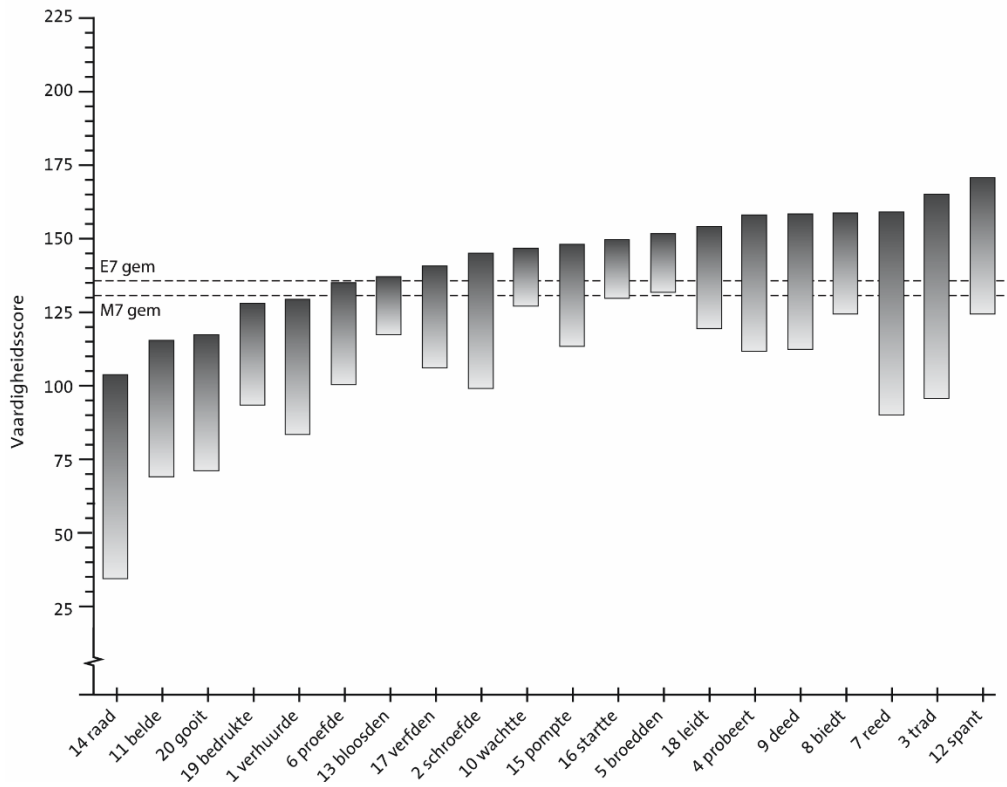
### Spelling niet-werkwoorden E7 digitaal - Taak 2



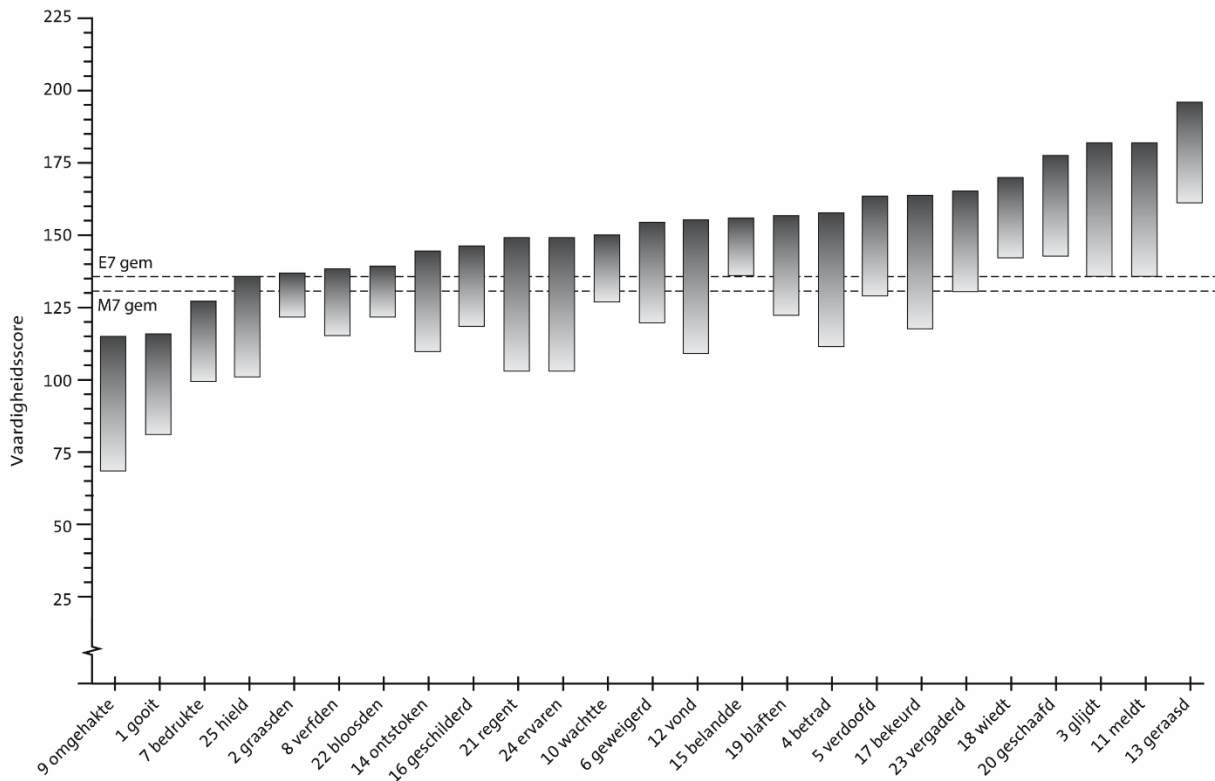
### Spelling werkwoorden M7 digitaal - Taak 1



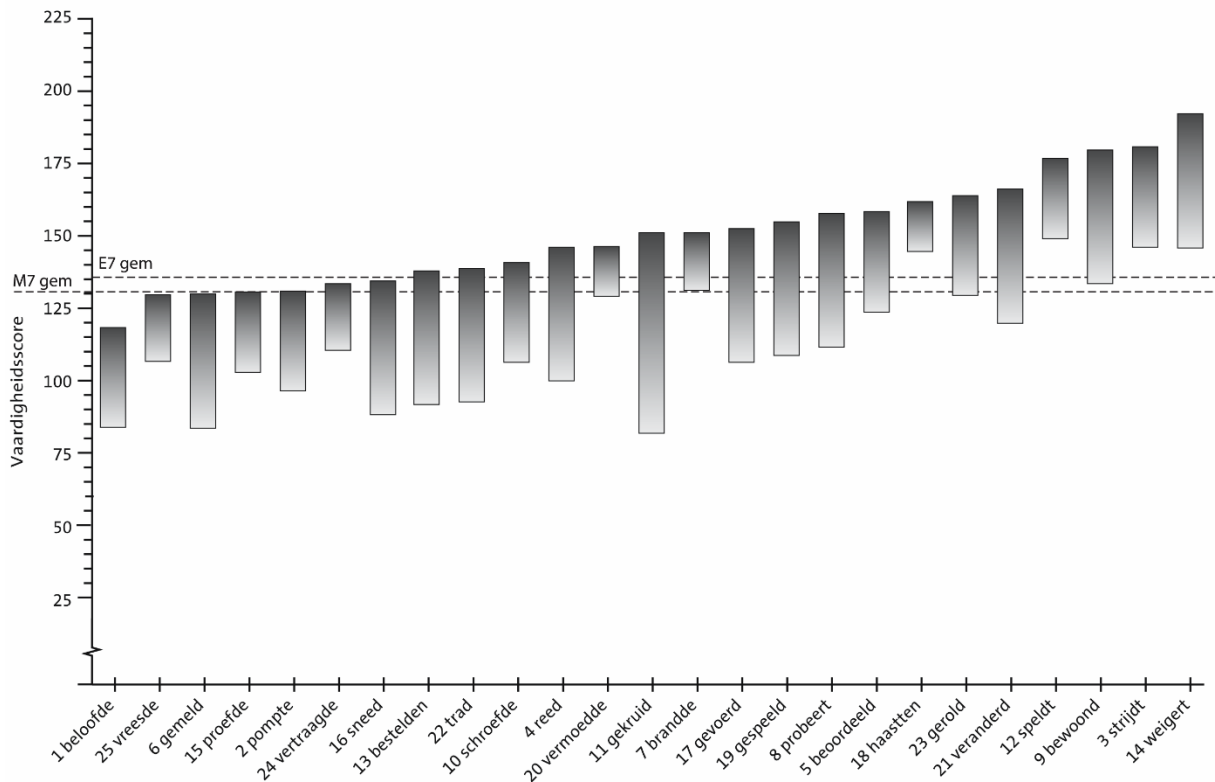
### Spelling werkwoorden M7 digitaal - Taak 2



### Spelling werkwoorden E7 digitaal - Taak 1



### Spelling werkwoorden E7 digitaal - Taak 2



## Bijlage 2 Klassieke en IRT-indices van de opgaven in digitale toetsen Spelling 3.0 groep 7

Toets M7 niet-werkwoorden

Volgnr	P-Val	RIT	Beta	Info
1	0,842	0,218	-0,544	0,506
2	0,802	0,333	-0,191	1,263
3	0,582	0,465	0,223	3,010
4	0,769	0,493	0,012	3,269
5	0,588	0,529	0,256	4,256
6	0,480	0,567	0,389	5,712
7	0,846	0,308	-0,308	1,054
8	0,598	0,529	0,245	4,229
9	0,515	0,381	0,338	1,902
10	0,835	0,392	-0,165	1,827
11	0,783	0,341	-0,146	1,345
12	0,507	0,524	0,348	4,369
13	0,654	0,576	0,189	5,269
14	0,577	0,465	0,231	3,018
15	0,698	0,451	0,085	2,660
16	0,601	0,464	0,191	2,972
17	0,534	0,382	0,308	1,896
18	0,845	0,385	-0,306	1,747
19	0,715	0,513	0,054	3,679
20	0,739	0,357	-0,052	1,507
21	0,802	0,333	-0,193	1,262
22	0,611	0,464	0,212	2,951
23	0,403	0,367	0,519	1,840
24	0,789	0,483	-0,021	3,092
25	0,830	0,396	-0,154	1,866
26	0,842	0,387	-0,181	1,770
27	0,760	0,350	-0,098	1,432
28	0,730	0,360	-0,036	1,537
29	0,886	0,348	-0,296	1,380
30	0,619	0,380	0,166	1,808
31	0,608	0,464	0,215	2,956
32	0,748	0,437	0,003	2,410
33	0,829	0,459	-0,092	2,691
34	0,646	0,526	0,185	4,055
35	0,843	0,448	-0,181	2,534
36	0,471	0,518	0,407	4,357
37	0,542	0,464	0,304	3,062
38	0,784	0,486	-0,012	3,139
39	0,552	0,464	0,291	3,052
40	0,659	0,375	0,096	1,729
41	0,703	0,450	0,023	2,635
42	0,588	0,465	0,243	2,999
43	0,846	0,308	-0,557	1,054
44	0,647	0,460	0,121	2,850
45	0,700	0,367	0,023	1,628
46	0,670	0,373	0,078	1,706
47	0,479	0,378	0,393	1,901
48	0,805	0,331	-0,397	1,248
49	0,758	0,433	-0,014	2,354
50	0,822	0,323	-0,239	1,174

Toets E7 niet-werkwoorden

<b>Volgnr</b>	<b>P-Val</b>	<b>RIT</b>	<b>Beta</b>	<b>Info</b>
1	0,881	0,291	-0,229	0,857
2	0,800	0,426	0,097	2,066
3	0,581	0,539	0,451	4,206
4	0,746	0,261	-0,042	0,705
5	0,590	0,390	0,402	1,832
6	0,684	0,465	0,292	2,686
7	0,572	0,391	0,430	1,852
8	0,638	0,472	0,358	2,841
9	0,783	0,351	0,040	1,333
10	0,566	0,474	0,457	2,997
11	0,606	0,283	0,312	0,877
12	0,817	0,531	0,154	3,657
13	0,641	0,537	0,376	4,013
14	0,496	0,468	0,551	3,042
15	0,769	0,440	0,154	2,264
16	0,674	0,382	0,257	1,681
17	0,742	0,568	0,265	4,503
18	0,599	0,474	0,412	2,939
19	0,764	0,442	0,162	2,291
20	0,680	0,466	0,298	2,698
21	0,596	0,474	0,403	2,946
22	0,713	0,459	0,247	2,559
23	0,673	0,467	0,258	2,726
24	0,774	0,355	0,061	1,370
25	0,896	0,408	-0,052	1,846
26	0,839	0,322	-0,101	1,081
27	0,803	0,539	0,174	3,834
28	0,529	0,535	0,511	4,289
29	0,612	0,389	0,365	1,802
30	0,772	0,439	0,148	2,247
31	0,550	0,391	0,467	1,870
32	0,709	0,460	0,254	2,578
33	0,564	0,474	0,463	3,000
34	0,859	0,309	-0,157	0,979
35	0,553	0,473	0,475	3,013
36	0,771	0,505	0,148	3,207
37	0,477	0,465	0,580	3,037
38	0,560	0,538	0,478	4,250
39	0,694	0,378	0,222	1,632
40	0,597	0,474	0,418	2,945
41	0,696	0,378	0,218	1,626
42	0,857	0,449	0,041	2,342
43	0,487	0,528	0,565	4,299
44	0,637	0,472	0,362	2,846
45	0,638	0,588	0,394	5,259
46	0,731	0,454	0,147	2,471
47	0,623	0,388	0,347	1,785
48	0,823	0,412	0,048	1,900
49	0,764	0,560	0,234	4,284
50	0,878	0,369	-0,083	1,448



Toets M7 werkwoorden

<b>Volgnr</b>	<b>P-Val</b>	<b>RIT</b>	<b>Beta</b>	<b>Info</b>
1	0,854	0,243	-0,358	0,733
2	0,807	0,268	-0,248	1,289
3	0,587	0,323	0,152	1,696
4	0,640	0,227	-0,023	0,874
5	0,843	0,176	-0,594	0,403
6	0,570	0,465	0,207	4,707
7	0,396	0,544	0,373	7,513
8	0,640	0,227	-0,019	0,775
9	0,706	0,303	-0,038	1,336
10	0,487	0,325	0,309	1,896
11	0,848	0,173	-0,612	0,389
12	0,446	0,512	0,334	5,830
13	0,662	0,451	0,117	3,231
14	0,555	0,466	0,223	4,739
15	0,685	0,220	-0,123	0,708
16	0,518	0,595	0,267	9,704
17	0,476	0,234	0,335	0,925
18	0,620	0,553	0,190	5,586
19	0,739	0,294	-0,097	1,213
20	0,791	0,195	-0,411	0,515
21	0,768	0,284	-0,154	1,101
22	0,683	0,309	-0,007	1,752
23	0,647	0,226	-0,038	0,867
24	0,605	0,321	0,121	1,919
25	0,453	0,555	0,321	7,967
26	0,714	0,375	0,014	2,042
27	0,672	0,222	-0,098	0,838
28	0,529	0,403	0,245	3,300
29	0,601	0,322	0,136	1,660
30	0,509	0,560	0,273	8,023
31	0,830	0,256	-0,304	1,175
32	0,521	0,326	0,247	1,997
33	0,618	0,554	0,177	7,652
34	0,856	0,170	-0,652	0,478
35	0,617	0,396	0,145	2,559
36	0,476	0,558	0,300	8,011
37	0,673	0,386	0,108	2,955
38	0,568	0,401	0,207	2,771
39	0,759	0,359	-0,061	2,503
40	0,822	0,260	-0,284	1,215

Toets E7 werkwoorden

<b>Volgnr</b>	<b>P-Val</b>	<b>RIT</b>	<b>Beta</b>	<b>Info</b>
1	0,865	0,227	-0,284	0,992
2	0,738	0,510	0,190	6,696
3	0,569	0,307	0,291	1,991
4	0,658	0,298	0,152	1,838
5	0,581	0,381	0,291	3,295
6	0,651	0,437	0,233	4,520
7	0,825	0,316	-0,061	2,049
8	0,753	0,349	0,108	2,589
9	0,875	0,220	-0,318	0,930
10	0,635	0,535	0,273	7,798
11	0,518	0,307	0,360	2,023
12	0,687	0,208	-0,019	0,822
13	0,390	0,426	0,504	4,701
14	0,731	0,356	0,095	2,717
15	0,529	0,496	0,361	6,551
16	0,729	0,419	0,142	3,997
17	0,634	0,302	0,184	1,890
18	0,531	0,308	0,346	2,019
19	0,580	0,306	0,269	1,978
20	0,464	0,305	0,432	2,016
21	0,728	0,200	0,033	0,760
22	0,737	0,510	0,177	6,717
23	0,567	0,382	0,308	3,319
24	0,721	0,286	0,040	1,657
25	0,790	0,264	-0,097	1,381
26	0,858	0,232	-0,262	1,035
27	0,702	0,364	0,145	2,876
28	0,488	0,380	0,403	3,372
29	0,717	0,203	-0,098	0,779
30	0,638	0,376	0,223	3,144
31	0,843	0,240	-0,222	1,116
32	0,563	0,497	0,334	6,481
33	0,674	0,296	0,121	1,797
34	0,534	0,382	0,349	3,360
35	0,746	0,279	-0,007	1,566
36	0,778	0,269	-0,070	1,434
37	0,393	0,427	0,498	4,715
38	0,766	0,273	-0,038	1,487
39	0,527	0,220	0,335	0,948
40	0,788	0,335	0,014	2,346
41	0,786	0,266	-0,085	1,402
42	0,732	0,356	0,096	2,713
43	0,428	0,523	0,442	8,165
44	0,690	0,293	0,094	1,755
45	0,652	0,568	0,267	9,320
46	0,624	0,303	0,202	1,910
47	0,693	0,207	-0,038	0,814
48	0,578	0,382	0,295	3,301
49	0,747	0,464	0,149	5,178
50	0,803	0,328	-0,019	2,229



**Cito**

Amsterdamseweg 13  
6814 CM Arnhem  
Postbus 1034  
6801 MG Arnhem  
T (026) 352 11 11

Fotografie: Gijs Versteeg

