



#pkdigitaal

Stavaza project: persoonskaarten overledenen 1939 – 1994 digitaliseren

Van karton naar data....



Agenda

1. **Deelresultaat digitalisering**
2. Deelresultaat zoekingang / interne database
3. Afsluitende opmerkingen

PERIODEN	AANTAL PERSOONSKAARTEN	METER
1939 – 1970	2.880.000	758
1971 – 1980	1.130.000	340
1981 – 1987	920.000	275
1988 – 1994 (tot 1 oktober)	870.000	260
TOTAAL	5.800.000	1.633

Verhouding standaard persoonskaart / kwetsbare persoonskaart

- Standaard kartonnen persoonskaart 99% = 5.740.000 (met een afwijking van +5/-5%)
- Kwetsbare persoonskaart ruim 1% = 60.000 (met een afwijking van +50/-50%)

Van der A-- 2 318 , Ghebt		Vader - wief Hendrick - wief	
Hendricus Pieter Nicolaas-- 25 Januari 1941 p. s-Gravenhage		JK Lisveruan Gravenhede	
Pieter Gerardus-- Wifer Johanna-- 3 Nov 17 = Ov 24 Apr 19 = 31			
Leliveld, Theodora Jacoba Maria Aikenaar, Magdalena Johanna		26 apr 45 = 1 nov 55 = 19 apr 82 = 8 Ov 27 Jan 48 = 24 sep 89 = Ov	
20dec01 Anbachtsegaard 114 14okt82 Drieff Footstraat 173 14jun84 Nouterpad 26 17okt84 Westinde 35 14jun89 Westendorp wuid 504 20dec89 Chopinlaan 152		Geboren in 1887 es 2/10294 562 4796 954 1911 1912 1913 1914 1915 1916 1917 1918 1919 1920 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932 1933 1934 1935 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972 1973 1974 1975 1976 1977 1978 1979 1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030 2031 2032 2033 2034 2035 2036 2037 2038 2039 2040 2041 2042 2043 2044 2045 2046 2047 2048 2049 2050 2051 2052 2053 2054 2055 2056 2057 2058 2059 2060 2061 2062 2063 2064 2065 2066 2067 2068 2069 2070 2071 2072 2073 2074 2075 2076 2077 2078 2079 2080 2081 2082 2083 2084 2085 2086 2087 2088 2089 2090 2091 2092 2093 2094 2095 2096 2097 2098 2099 2100	
Wan42 Mariastraat 2a Apr42 Loostrat 125 Mar48 Koninginnegrat 23 19okt48 Loosduinseweg 523 19okt48 Loosduinseweg 571 10jul59 Groenteweg D3 (dec57) 19sep59 van Kerveldstrat 24 Apr63 Lepentenselaan 258 6me164 Fahrenbelstraat 5 14okt66 Hendrik Mandestraat 14okt66 Van der Meerstrat 11 19sep66 Stadelde 8 8dec66 Jacoba Mariastraat 78 14jun75 BOOGERMAN HATZELDORP Willems 3 20dec74 81 Lentenselaan 202			

Van der A, Roland		10 aug 67	
Joni		21 aug 72	
Schalk, Leo Eugenia		15 sep 69	
Alexander		3 Feb 72	

Gerardus Bernardus-- 20 September 1901		JK	
Sijkes, Antonia Hendrika-- 10 Aug 90 = Arden 3 Sep 98 = Arden			
Eenr. arbeider			
ANNO 1906 13okt90 Janskerstrat 7 13okt90 Janskerstrat 7 27aug11 JB 448/614-c			

21 aug 67		21 aug 72	
15 sep 69		3 Feb 72	

Standaard en kwetsbare
persoonskaarten digitaliseren



#pkdigitaal

Het archief voordat team scanvoorbereidingen aan de slag gaat....

- Iedere persoonskaart wordt gecontroleerd: kwetsbare kaarten worden apart gelegd
- Witte kaarten worden verwijderd
- Plakband, nietjes, etc. worden verwijderd, op vervolgkaart wordt (indien nodig informatie toegevoegd met elek. typemachine)



Het archief nadat team scanvoorbereidingen haar werk heeft gedaan....

- Natuurlijk proberen we de negatieve invloed op de dienstverlening voor de klant zo veel mogelijk te beperken.



Het archief nadat de dozen weer terugkomen na digitaliseren....

- Iedere doos heeft een etiket met een uniek doos-id en barcode
- De eerste PK van iedere doos wordt gemetadateerd in een google sheet. (autorisatie via Jeroen)
- We gaan vanaf 2/9 naar twee transporten per week (donderdagochtend en donderdagmiddag) wekelijks 200-300 dozen

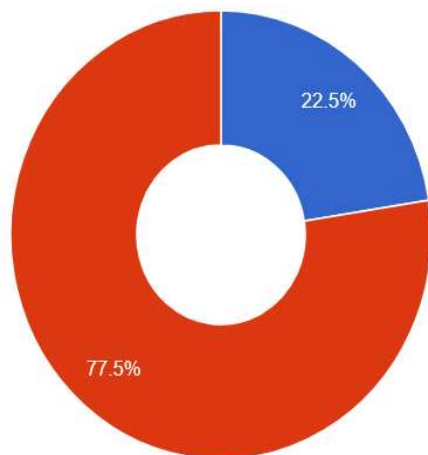


2 september vieren we de mijlpaal:

1 miljoen goedgekeurde persoonskaarten digitaal..



Algehele project status



● Inmiddels gescand ● Nog te doen

Totaal aantal kaarten gescand

1.306.936

Dozen

Aantal dozen opgehaald : **1615**

Aantal dozen retour gebracht : **1465**

Aantal dozen bij MultiScan : **150**

Productieschema

26-aug	14	150	830	124.500	1.351.110	
2-sep	15	200	830	166.000	1.517.110	tweemaal
9-sep	16	200	830	166.000	1.683.110	
16-sep	17	250	830	207.500	1.890.610	
23-sep	18	250	830	207.500	2.098.110	
30-sep	19	250	830	207.500	2.305.610	
7-okt	20	250	830	207.500	2.513.110	
14-okt	21	250	830	207.500	2.720.610	
21-okt	22	300	830	249.000	2.969.610	
28-okt	23	300	830	249.000	3.218.610	
4-nov	24	300	830	249.000	3.467.610	
11-nov	25	300	830	249.000	3.716.610	
18-nov	26	300	830	249.000	3.965.610	
25-nov	27	300	830	249.000	4.214.610	
2-dec	28	300	830	249.000	4.463.610	
9-dec	29	300	830	249.000	4.712.610	
16-dec	30	300	830	249.000	4.961.610	
23-dec	31	150	830	124.500	5.086.110	
6-jan	32	200	830	166.000	5.252.110	
13-jan	33	300	830	249.000	5.501.110	
20-jan	34	300	830	249.000	5.750.110	
27-jan	35	120	500	60.000	5.810.110	Kwetsbaar
		7.035		5.810.110		





PK beeldbestanden duurzaam bewaren en opnemen in huidige dienstverlening

Masterbestanden duurzaam offline bewaren in een Microsoft Azure datacenter.

Afgeleide bestanden beschikbaar stellen voor huidige dienstverlening.

- Sharepoint omgeving (vervallen)
- Digitaal Depot Multiscan. Het d-depot is een digitale opslagfaciliteit van MultiScan, waarbij op de servers van MultiScan ruimte beschikbaar wordt gesteld aan het CBG. Geautoriseerde medewerkers krijgen via een beveiligde verbinding toegang tot het d-depot door middel van een unieke gebruikersnaam en wachtwoord.
- Censura (door Multiscan ontwikkeld) is een softwareprogramma dat is ontwikkeld voor het veilig, efficiënt en nauwkeurig anonimiseren van documenten.



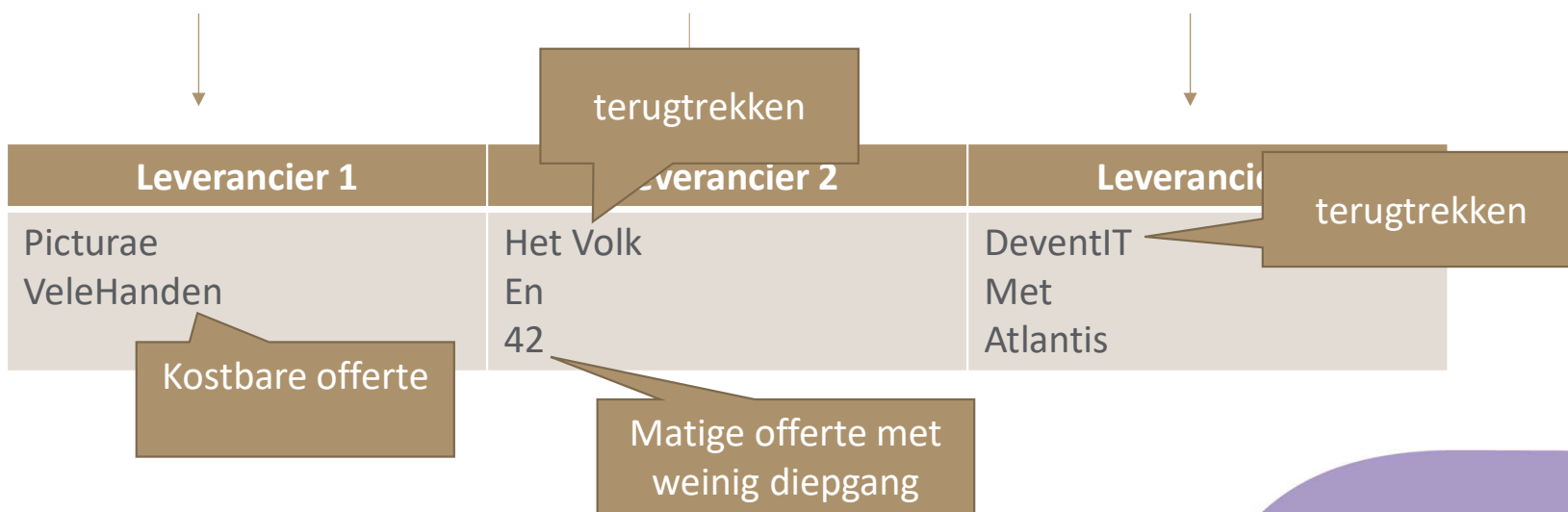
Agenda

1. Deelresultaat digitalisering
2. **Deelresultaat zoekingang / interne database**
3. Afsluitende opmerkingen

PvE en offertetraject

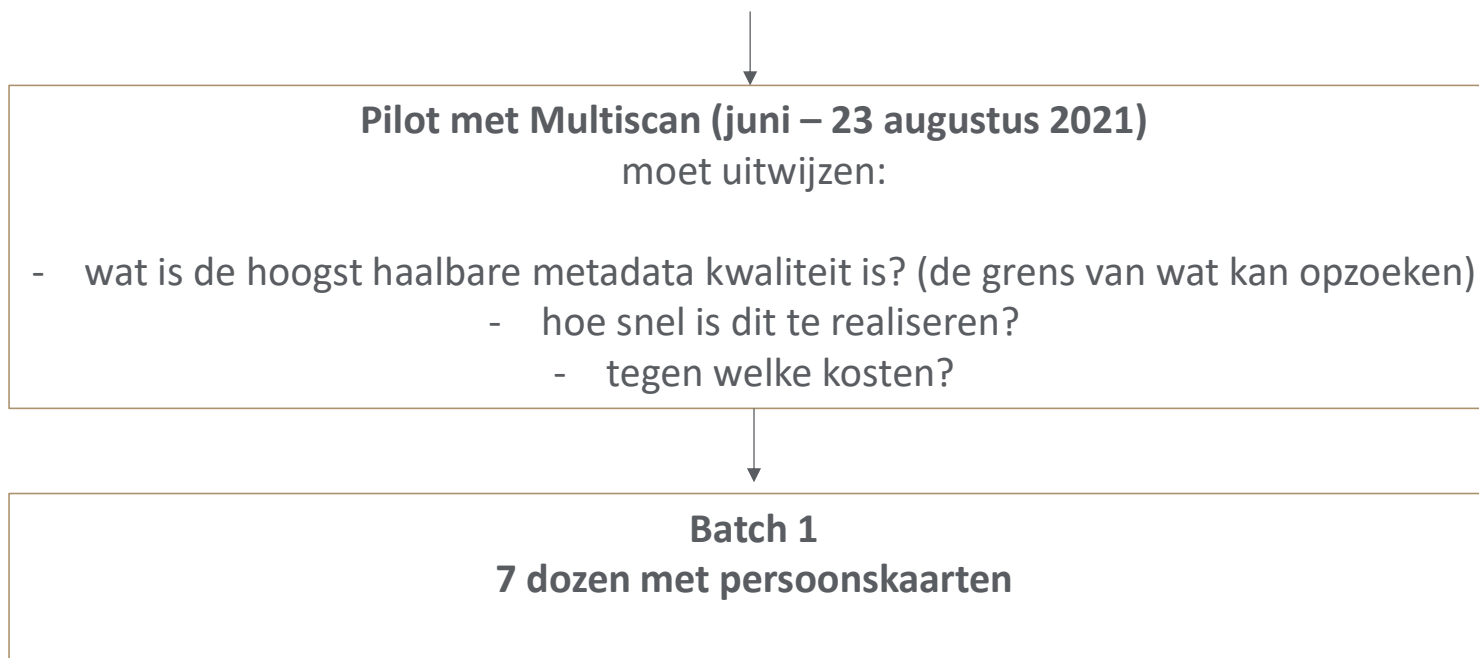


PvE Crowdsourcing traject en tijdelijk zoekstelsel
 Februari / maart 2021
 Offertetrajecten april / mei 2021



Nader onderzoek mogelijkheden en beperkingen slim OCR/HTR-en

Hoe kunnen we digitale persoonskaarten zo snel mogelijk, zo veel mogelijk geautomatiseerd, transcriberen naar betrouwbare metadata?



Nader onderzoek mogelijkheden en beperkingen slim OCR/HTR-en

Hoe kunnen we digitale persoonskaarten zo snel mogelijk, zo veel mogelijk geautomatiseerd, transcriberen naar betrouwbare metadata?



Multiscan heeft een partner Scalehub in de arm genomen voor deze pilot:

begin 2016 opgericht
een internationaal multicultureel bedrijf, actief over de hele wereld
uitgebreide ervaring op het gebied van data automatisering, bedrijfsprocesbeheer en crowdsource technologieën, is ingezet om de traditionele data-extractie en documentverwerking revolutionair te veranderen

combineert het beste van menselijke vaardigheden en vindingrijkheid met bewezen automatiseringstools.
gebruikt de meest up-to-date methoden voor het converteren van beeldbestanden naar bruikbare gegevens met hoge snelheid en nauwkeurigheid
scherpe gegevensextractie en expertise voor het verwerken van grote volumes



Analyse en prioritering metadata Persoonskaart

Prio 1	Prio 2	Prio 3
Familienaam	Beroep	Bijzondere vermeldingen
Voornamen	Geboortedatum echtgeno(o)t(e)	Religie
Geboortedatum	Geboorteplaats echtgeno(o)t(e)	Aantekeningen
Geboorteplaats	Achternaam kinderen	Datum en plaats huwelijk gesloten (idem ontbinding)
Familienaamechtgeno(o)t(e)	Voornamen kinderen	A nummer (voorloperbsn nummer)
Voornamen echtgeno(o)t(e)	Geboortedatum kinderen	Registratie woonadressen
Overlijdensdatum	Geboorteplaats kinderen	Etc.
Overlijdensplaats	Relatie van het kind tot het gezinshoofd, bijvoorbeeldz (zoon), d (dochter), sz (stiefzoon) ofsd (stiefdochter)	
Periode	Achternaam en voornamen ouders	

Scope pilot: focus prio 1, maar ook prio 2 en 3 i.v.m. onderkennen laaghangend fruit

Dia 16

JT1

Jan Thielen; 31-8-2021

Pilotresultaten: automatisch herkennen typen persoonskaarten

Model kaart	Aantal kaarten	Percentage
Persoonskaarten	5.288	92,5%
<i>Model A</i>	4.603	80,5%
<i>Model B</i>	685	12%
<i>Verwijskaart</i>	129	2,3%
Onbekend*	297	5,2%
Totaal	5.714	100%

*Kaarten zijn niet automatisch geclassificeerd. Dit zal handmatig uitgevoerd worden.

Nader onderzoeken

Pilotresultaten: OCR kwaliteit in %

Prioriteit 1

Veld	OCR-kwaliteit in %
3a – Familiennaam	99,00
3b – Voornamen	92,00
4 – Geboortedatum	89,00
4 - Geboorteplaats	97,00
9 – Familiennaam echtgeno(o)t(e)	89,00
9 – Voornamen echtgeno(o)t(e)	89,00
Overlijdensdatum*	84,00
Overlijdensplaats*	92,00

*Kwaliteit gebaseerd op die kaarten waarin deze gegevens zijn herkend.

*Percentage van kaarten waarop deze gegevens zijn herkend. Van 5775 zijn er 2317 kaarten waar geen stempel is herkend. Dit betekent niet dat er geen stempel aanwezig is.

Prioriteit 2

Veld	OCR-kwaliteit
7 - Beroep	99,60
11 - Geboortedatum echtgeno(o)t(e)	79,80
11 -Geboorteplaats echtgeno(o)t(e)	99,80
28 - Achternaam kinderen	80,50
28 - Voornamen kinderen	81,00
30 - Geboortedatum kinderen	96,00
30 - Geboorteplaats kinderen	89,00
32 - Relatie van het kind tot het gezinshoofd	98,00
8 - Achternaam en voornamen ouders	99,85
A-nummer (voorloper BSN-nummer)	Niet bekend

Prioriteit 3

Veld	OCR-kwaliteit
23,24,25, 35 - Bijzondere vermeldingen	Niet te definiëren, te divers om hier een percentage aan te hangen
6 - Religie	99,50
33,34 - Aantekeningen	Niet te definiëren, te divers om hier een percentage aan te hangen
13, 14 - Datum en plaats huwelijk gesloten/ ontbinding	96,00
22 - Registratie woonadressen	70,00

In september 2021 controleert het CBG de pilotkwaliteit van de metadata



Pilotresultaten: wat is kwaliteit?

ABSOLUUT

Aantal karakters goed vertaald,
uitgedrukt in % van totaal aantal
karakters

Voor de prio 1 velden is het streven zo
dicht mogelijk bij 100% kwaliteit.
Streven naar 99,9%?
(nog af te spreken)

RELATIEF

Kwaliteit die goed genoeg is voor het
herkennen van de juiste persoonskaart:

3a - Familienaam	Herkende familienaam
Aalfs	aAalfs
Aalderink	X Aalderink
Aarntzen	A_Aarntzen



Agenda

1. Deelresultaat digitalisering
2. Deelresultaat zoekingang / interne database
3. **Afsluitende opmerkingen**

Kortom

De pilotresultaten geven veel stof tot nadenken. De maand september 2021 gebruiken we voor het analyseren van de pilotresultaten en de vertaling naar een best haalbare planning die tevens rekening houdt met de beschikbare middelen.



Van ambitie naar gefaseerde realisatie

#pkdigitaal

Fase 1: planning en ambities

Digitaliseren: 5,8 miljoen PK's in 8 maanden digitaliseren (januari 2022)

Masterbestanden duurzaam offline bewaren in een Azure netwerk

Afgeleide (JPEG) bestanden online beschikbaar voor verbetering huidige dienstverlening (efficiënter **richting kostendekkend** en hogere beeldkwaliteit voor de klant)

OCR / HTR 5-10% metadata (prio 1) + eerste versie tijdelijk zoekstelsel (maart 2022)

Fase 2: ambities

OCR / HTR stapsgewijs naar 100% metadata (10%, 20%, 33%, 50%, 100% prio 1) + verbeterd tijdelijk zoekstelsel (december 2022)

Experimenteren met automatische matching (juli – december 2022)

Fase 3: ambities

Toekomstvast NRO inclusief automatische matching (voorjaar 2023)

Keuzes maken: prio 2 en 3 velden metadateren

Dekking dankzij projectbijdrage vanuit Ministerie van Binnenlandse Zaken

Ambities vertalen naar planning, nieuwe subsidieaanvraag



Onzekere factor: corona invloed