# IMPACT: Centre of Competence in digitisation

*Hildelies Balk – Pennington de Jongh, IMPACT Project director, KB National Library of the Netherlands*
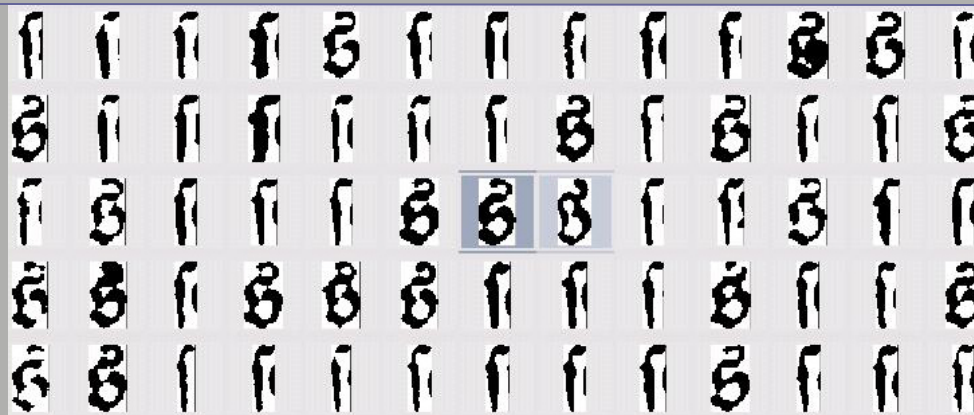
Improving Access to Text

**IMPACT**

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

# A short introduction...

http://www.youtube.com/user/theimpactproject

# Historical text: typical OCR results



**VVt Venetien den 1.Junij, Anno 1618.**

DJgn i f paffato te S' aö'Jifeert mo?üen/bah .)

etgi'uotbciraetail)i.r/JtmelchontDecht

te /sbnbe bele btr felbrr geiufttceert baer bnber

eeniglje jprant o^fen/bie ftcb .met beSpaenfcbeu

enbeeemgljen bifet Cbeiiupcen berbonbru befe

## OCR Challenges: damaged pages, bleed through,difficult layout, historic fonts … and many more

Improving Access to Text
# IMPACT

 IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

# Language Challenges: Spelling variants, orthographical variants, inflected forms…and more



## Historical variants of the Dutch word 'wereld' (world):

werelt weerelt wereld weerelds wereldt werelden weereld werrelts waerelds weerlyt wereldts vveerelts waereld weerelden waerelden weerlt werlt werelds sweerels zwerlys swarels swerelts werelts swerrels weirelts tsweerelds werret vverelt werlts werrelt worreld werlden wareld weirelt weireld waerelt werreld werld vvereld weerelts werlde tswerels werreldts weereldt wereldje waereldje weurlt wald weëled

Improving Access to Text

**IMPACT**

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

## Institutional Challenge: lack of knowledge and expertise → inefficiency

Improving Access to Text

**¡MPACT**

KB  IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

# Answering the challenges – IMPACT

- IMPACT – Improving Access to Text (2008-2011)

  Large-scale integrating research project, funded by the EC
  - Consortium of 26 partners
    - Good mix of public and private partners
    - Users, researchers and industry work together to find solutions
    - Each established in a large international network
  - Coordinated by the National Library of the Netherlands (KB)
  - EU funding: € 12 100 000 (FP7 ICT Work Programme)
  - From 2012: sustainable Centre of Competence with alternative resources

- Main objectives:

  - Innovate OCR & language technology

  - Build capacity in mass-digitisation

# IMPACT Main Achievements at this point

- Improved commercial OCR (ABBYY 'IMPACT' Finereader 10 on market)
- Effective tool for OCR correction with volunteer involvement (CONCERT) ready for implementation
- Novel Approaches to preprocessing, OCR and post correction available
- Computerlexica for nine languages close to delivery
- Digitisation Framework with evaluation tools available
- Facility to plug in other tools (Your tools!)
- Large Dataset with sophisticated Ground Truth close to final delivery
- Knowledge bank with guidelines and learning resources close to delivery
- A unique network bringing together experts from different communities
- New website under construction
- Centre of Competence to be launched at final conference 24-25 October 2011

# A Centre of Competence in digitisation: Why

- Challenges in digitisation of historic material still there
- No lack of novel approaches to improve access: IMPACT and many others
- Challenge: implementation in real life
- Challenge: direction of research
- Need for real life datasets with Ground Truth
- Need for real life testing and evaluation
- Need for support in implementing the  solutions

Improving Access to Text

# IMPACT

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

# Building the Business Model: Methodology

Business Model Generation (http://www.businessmodelgeneration.com)

| Key Partnerships | Key Activities | Value Propositions | Customer Relationships | Customer Segments |
|---|---|---|---|---|
| | Key Resources | | Channels | |
| Cost Structure | | Revenue Streams | | |

Improving Access to Text

# IMPACT

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

## Building the Business Model: sessions

# Building the Business Model: Gathering ideas

# Building the Business Model: Finding commitment

## CONTENT HOLDERS

Blue blocks: shared with all customers segments

Red blocks: focus on Content Holders

## Key Activities

Coordination of all activities in the Centre

Coordination of requests that come in through website

Maintaning the framework

Maintaning the evaluation tools

Maintain Datasets In libraries and at the providng organisation

Finding new funding opportunities and engaging in new projects

Website Content Management

Experts serving requests that come in through website

Providing and maintaining Online Tutorials

Maintaining knowledge base

Marketing in MLA community at different levels by taking part in international forums and conferences, and bilateral contacts

Engaging with volunteers for coorection of OCR

Providing in house training courses on request

## Key Partnerships

Training Providers

Content holders around core of IMPACT libraries

Organisations in MLA field

Funders

Consortia of related projects

Partners with business & marketing expertise from IMPACT consortium

Content Providers (e g e-books mobile)

Educators (digital Library based courses)

Search Engine companies

## Key Resources

A very small permanent CoC team: one full itme expert intermediary for the network and one part time admin

Human Experts (affiliated to centre/IMPACT partner institutions)

Office with communication Infrastructure

Helpdesk

Website

Datasets with metadata and ground truth

Evaluation tools

Framework

## Value Propositions

Free (registration)

One stop shop for digitisation: information sharing, tool demos, online training

Community Network of experts in digitisation

Access to network of experts in research

Research Dataset access to images demoset GT

Access to language resources

Online Tutorials

Knowledge base

Taverna Framework as is with a good user manual

Subscription fee of € 000- 10 000 euros a year (dependent of budget)

Membership of CoCboard – opportunity to shape Centre, shape discussions on stadards, decide on new partners

Opportunities to send (young) staff members abroad for short periods to help set up digiprojects

Training & Education in digitisation on CoC membership conditions (e.g. one course a year free)

Evaluation Service (third party tools & production workflows)

Dataset resarch & evaluation (access to images Ground Truth OCR and metadata Community forum benchmarking)

Reduced fee for pay as you go services that are not included

Framework Integration in your library workflow with technical support from CoC (utilising Taverna toolkit)

Access to technical support from research partners (on CoC SLA basis)

Certain level of on demand tool services (e g small scale OCR)

Pay as you go

On demand tool services for non members (e g small scale OCR and enhancement Evaluation toolkit)

in person Training & Education in digitisation

Digitisation consultancy by experts from partners (government /local funding)

## Customer Relationships

Centralized in CoC: One Expert technical/ executive cooridnator handles requests for service, giving (technical) advice, maintaining customer relatinships and acquiring new customers

Contact Point (HelpDesk) for first basic information

Expert to Expert (research) Expert to expert (digitisation)

Dedicated Contact / Consultant

Committees & Working Groups (e.g. standards) to be set up and maintained by interested partners

## Channels

Website

Publications by partners

Conferences

TEL/ Europeana

Direct Contacts

Library organisations (CENL, LIBER, IFLA, CERL)

Mailing Lists

## Customer Segment: Content Holders

Private sector content holder (e.g. Publishers)

Public sector content holders (e.g. MLA, Scientific)

National Libraries as institutions

Researchers interested in content (historians, language experts, journalists, amateur genealogists) – could pay for small scale OCR and enhancement of text

## Cost Structure

Costs for (technical) Coordination and services 1 fte

Costs for at least 0 5 fte admin support

Hosting datasets (one partner) Server (+software)

Hosting and maintaining framework (one partner) infrasturcture, hardware0, 2 fte

Maintaning datasets hosting partner 0,2 fte libraries 0,1 fte each

Hosting and maintaining website (one or two partners) Infrastructure, hardware, (supporting software?) and 0,2 fte

Travel and sustanance for digitisation experts

Office Costs

## Revenue Streams

Governments/ Public Funds for building digi Capacity

Pay-as-you-go usage of services

Membership Fees 10.000 ayer for board members

Fees for conferences and wrokshops

Contributions in Kind:making digitisation experts available for consultation and for updating and maintaining knowledge base and tutorials

Contributions in Kind:hosting the CoC office

Training / Accreditation fees

Funding (Partners / Public)

Contributions in Kind: volunteers correcting OCR

Contributions in Kind maintaining datasets (libraries)

Contributions in Kind: hosting webiste

15

**IMPACT**

# Centre of Competence: What*

- Not for profit organisation

- Mission: to support the process of making Europe's heritage accessible in digital form

- Web based international community with small core facility for support

- Curates IMPACT achievements and provides tools, services and facilities for further advancement of the State of the Art in this field

- Focuses on practical solutions

- Distributed effort by IMPACT partners spreads risk and ensures continuing engagement

- Income generation by offering number of resources and services at a fee (mix of subscription and pay as you go)

*pending approval of IMPACT consortium

# impact
### digitisation.eu
### centre of competence

About    Knowledge bank    Tools    Services    News & Events    Blog    Community

## Welcome to **impact**
## Centre of Competence

Lorem ipsum dolor sit amet, consecte adipi scing elit. Etiam eleifend quis diam tincidu sagittis. Curabitur ut augue massa.

*Lorem ipsum dolor sit amet, consecte adipi scing elit etiam*

**TWITTER** : Cras tristique massa vel metus ornare ut vulputate eros pellentesque. Integer non magna est. Integer mi urna, congue sed adipiscing vitae

## What's new at **impact**

Sapien velit interdum dui, non lobortis dolor leoquis odio. Curabitur nec lacus etx eget nisi tempus vehicula **more >**

Sapien velit interdum dui, non lobortis dolor leoquis odio. Curabitur nec lacus etx eget nisi tempus vehicula **more >**

Sapien velit interdum dui, non lobortis dolor leoquis odio. Curabitur nec lacus etx eget nisi tempus vehicula **more >**

Sapien velit interdum dui, non lobortis dolor leoquis odio. Curabitur nec lacus etx eget nisi tempus vehicula **more >**

### What do **impact** do?

**Click here to watch our film**

### Tools

Lorem ipsum dolor sit amet, consecte adipi eleind sagittis. Curabitur ut augue massa. Quue **more >**

### Knowledge Bank

Lorem ipsum dolor sit amet, consecte adiifend sagittis. Curabitur ut augue massa. Quue **more >**

### Services

Lorem ipsum dolor sit amet, consecte adipifend sagittis. Curabitur ut augue massa. Quue **more >**

Improving Access to Text

iMPACT

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

# Centre of Competence: Who

- Digitisation practitioners in content holding institutions
- Researchers in the field of historical document processing and language technology
- Service providers
- End users (Researchers in the humanities who need access to the content)

- IMPACT Centre of competence offers distinct value to each of these target groups

# Centre of Competence*: benefits for content holders**

- Exchange of best practice in community of content holders with digitisation programmes

- KnowledgeBank with comprehensive and up to date information and technology watch reports

- Training on demand and on line tutorials

- Online support through a Helpdesk

- Support in the implementation of the innovative IMPACT solutions for improving access to text

- Access to the IMPACT Dataset with Ground Truth and tools for evaluation

- Digitisation Framework: Guidelines on using the open source workflow management system Taverna in a digitisation workflow

- Language Resources: a set of historical lexica that can be utilised within the digitisation workflow

- Conferences/workshops with focus on demonstration and implementation

- Working together on future practical solutions with scientific communities in the area of pattern recognition, language technology, image processing

*pending approval of IMPACT consortium  **access depends on membership level

# Centre of Competence*: Benefits for researchers**

- New community: Bridging the gap between specialist research and real life needs
- Brings together scientific communities in the area of pattern recognition, language technology, image processing with a focus on large scale digitisation
- Access to content holding community
- Access to large real life datasets and ground truth
- Working on implementation of research prototypes/products into digitisation environment
- Facilities for testing and evaluating new tools and IMPACT tools
- Working groups and comittees
- Access to new projects and funding opportunities
- Conferences/workshops with focus on demonstration and implementation

*pending approval of IMPACT consortium ** access depends on membership level

Improving Access to Text

**IMPACT**

IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

# Centre of Competence*: Benefits for service providers**

- Access to content holding community with large scale digitisation programmes
- Access to large real life datasets and ground truth
- Facilities for testing and evaluating new tools and IMPACT tools
- Attend conferences
- Working groups and comittees
- Make yourself known to your clients through yellow pages on website
- Sponsorsship and exhibition opportunities
- Working together with content holders and researchers on practical solutions

*pending approval of IMPACT consortium **access depends on membership level

# Join the Centre by becoming a member!

Three levels of membership (pending final decision IMPACT consortium) :

- Open (registration) access to forum, part of content
- Basic membership (fee): access to all facilities, reduced fee for conferences
- Premium membership (fee): member of the Board, privileges such as free entry to conferences

**Want to sign up?**

→ Mail to impact@kb.nl for information on membership

→ Join us now already on LinkedIn

→ Follow us on Twitter (@impactocr)

→ Access through www.impact-project.eu

→ Come to our conference on 24 and 25 October

# IMPACT final conference: 24-25 October 2011

Digitisation & OCR: Better, faster, cheaper
*Solutions of the IMPACT Centre of Competence and future challenges*

- Presentation of final results of IMPACT & related research in the area of OCR, digitisation and language technology

→ Location: The British Library, London, UK
→ Registration and more news available through the IMPACT website:
www.impact-project.eu

**Improving Access to Text**

# IMPACT

# www.impact-project.eu

**Improving Access to Text**

## IMPACT

printable view

"Search the IMPACT website"

Home
News
Helpdesk
Tools and applications
Calendar of events
About the project
Documents
Sitemap
Disclaimer
Contact
For partners

Twitter: @impactocr,
#impactproject

twitter  Linked in

WordPress  slideshare

You Tube  vimeo

IMPACT is a project funded by the European Commission. It aims to significantly improve access to historical text and to take away the barriers that stand in the way of the mass digitisation of the European cultural heritage. Read more

Friday 19. November 2010
**IMPACT at Czech Library conference**
In the beginning of the December, the IMPACT project will be promoted at the Czech national...

[more]