# Digidaily

## A newspaper dream come true

# The Collections at KB

- Swedish newspaper publishers deliver three legal deposit copies of all Swedish newspapers printed

- KB has 31 600 meters of newspapers or approx. 122 million pages

- 70 million pages are filmed on microfilm

- A major newspaper donation from a former journalist school
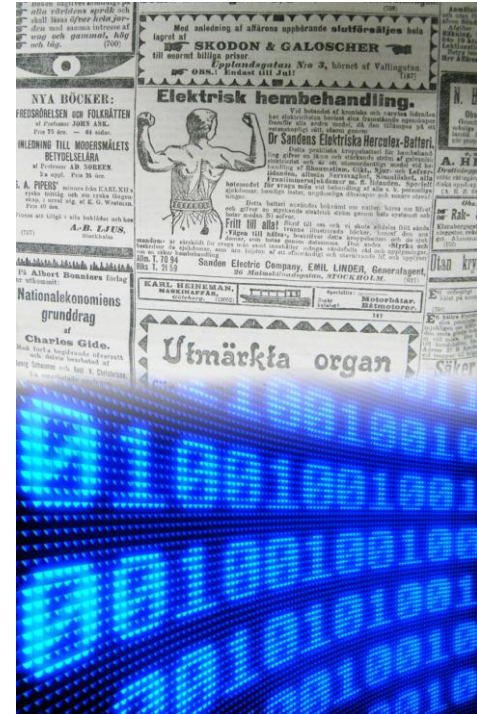
# The newspapers

# Activities at KB

- Requirement specifications

- Planning

- Preparation, registration and delivery

- Handling of rejected and damaged newspaper material

And much more…

# Requirement specifications

- **End product?**
  - What needs do the end-users have?
  - What metadata is needed?
  - What file formats should be used?
  = many different requirement specifications

- **Infrastructure**
  - How should the files be transferred from MKC to KB?
  - Storage at KB?

# Planning

- **Overall Planning**
  KB assesses the newspapers' condition,
  volume, size, font type, if it's already on
  microfilm, legal rights, scientific interest
  etc

- **Condition Survey**
  A condition check is done to get an
  overall picture of the newspapers'
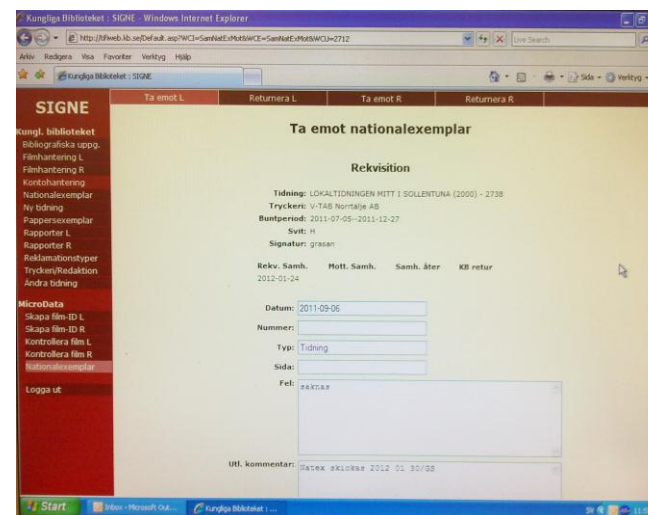  condition

# Preparation

- **Review of the material**
  - Cleaning
  - Preparation on issue level
  - Comments on condition, supplements, parts, editions, missing issues
  - If necessary when the material is damaged, find a replacement copy

# Registration

- Information on each bundle is recorded in KB's database

- Data export/import including newspapers name, date, comments about supplements, editions, condition etc

- Lists of expected editions and attachments are exported/imported

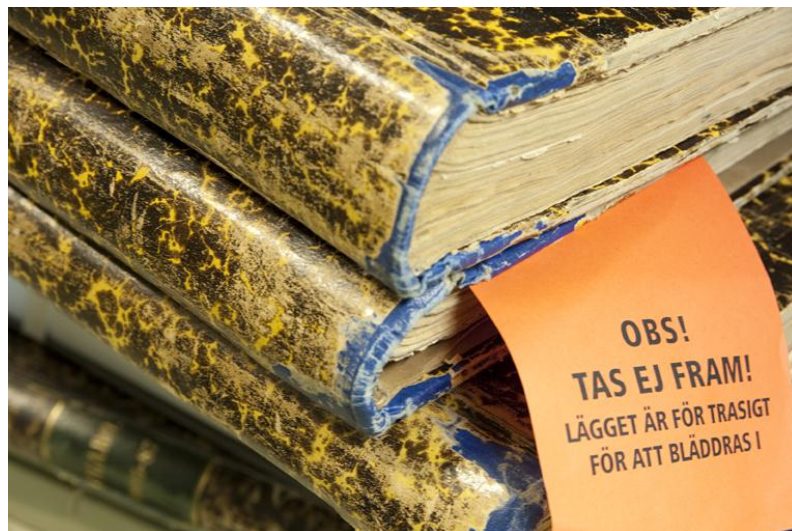- KB registers each bundle in the Workflow system

# Delivery

- For transportation of the official national copies special transport boxes has been built.

- Each box is followed by a printed packing list from the Workflow system.



**Digidaily**

# Damaged material

- **Rejects**
  What's the level for rejecting a page?
  Big holes? Small holes?
  Only loss of text?

- **Gaps**
  What gaps can be accepted?
  What's the priority? A complete
  material or large volume digitisation?

  **The goal is to find a rational way to
  handle rejects that suits both KB
  and MKC!**



OBS!
TAS EJ FRAM!
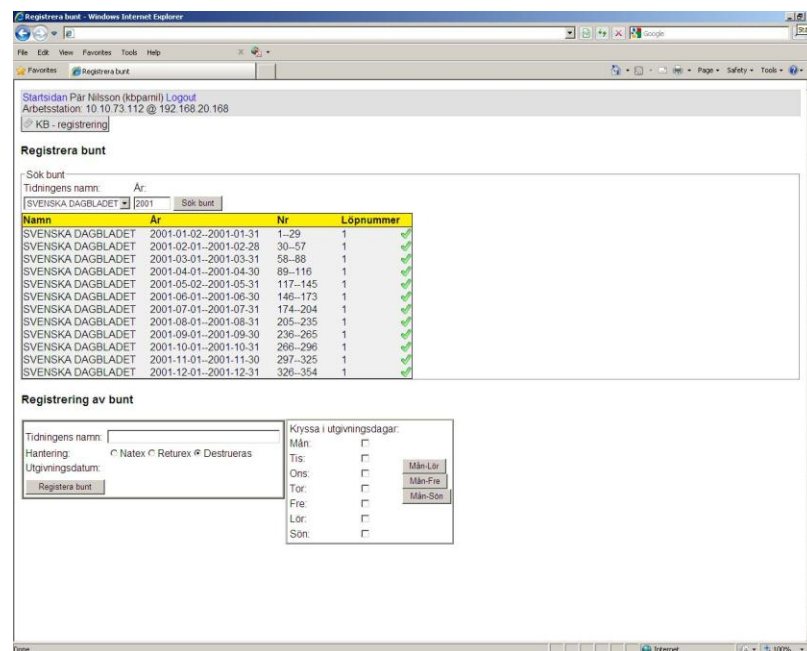LÄGGET ÄR FÖR TRASIGT
FÖR ATT BLÄDDRAS I

# Activities at MKC

- Develop workflow system

- Tests, evaluation and procurement (scanners, OCR etc.)

- Production planning

- Preparation and registration

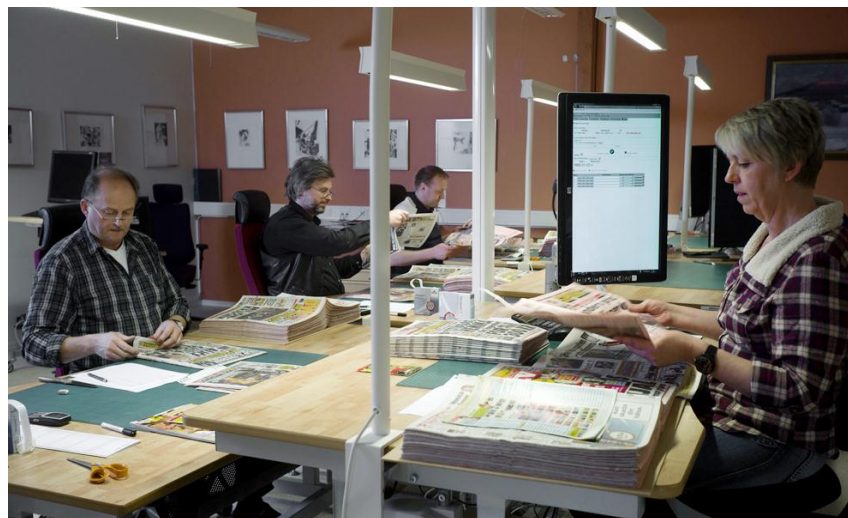- Scanning, OCR, control features

And much more…

# Workflow system

- Developed by a local team of developers at MKC

- Four major functions: database, add metadata through the process, keep track of and initiate the next process, collect data for statistics

- The workflow system is the backbone of the project

- Used by both MKC and KB

# Preparering på MKC

- The preparation process: go through and take apart between 650 and 6000 pages per person and shift

- The workflow system is used to collect metadata

Digidaily

# Scanning



SUPAG Mediascan 880c

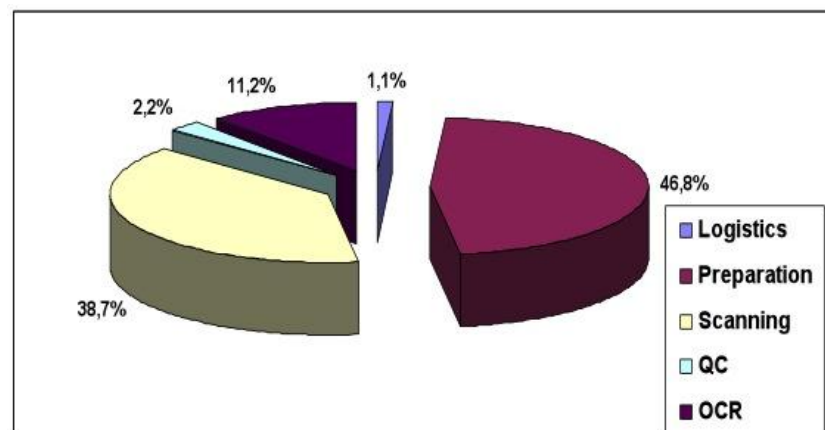Zeutschel OS

# OCR
## Optical Character Recognition

- Content conversion software by Zissor

- Segmentation on article level

- Word lists used for OCR to improve word accuracy.

- No manual correction on headlines or text

- OCRed text and metadata in ALTO

# £1 a page, one € a page or…

- Category 1 - bound, torn, fragile paper, large formats

- Category 2 - bound, where most can be taken apart and only a few are kept still bound, fair paper quality

- Category 3 – tabloids, stapled but not bound

- The price span today is between € 0.25–1.11/page incl. OCR

# Digidaily

– ett utvecklingsprojekt där
Riksarkivet, Kungliga biblioteket och Mittuniversitetet tillsammans ska
utveckla rationella metoder och processer för digitalisering av dagstidningar.
Projekttid 2010-04-01--2013-03-31

Finansiärer är EU:s strukturfonder i Mellersta Norrland, Riksarkivet, Kungliga biblioteket,
Mittuniversitetet, Schibsted Sverige och Länsstyrelsen i Västernorrland