

# DIGITALISERING AF 32 MIO. AVISSIDER PÅ 3 ÅR

Aalborg 19-11-2014

# OVERBLIK

- Kort om Statens Avissamling
- Avisdigitalisering – baggrund
- Avisdigitalisering – hvordan
- Tilgængelighed – hvordan får man adgang?
- Spørgsmål?

# KORT OM STATENS AVISSAMLING

# STATENS AVISSAMLING



*Statens Avissamling flytter ind. Demokratens fotograf var så heldig at være til stede, netop da ældre årgange af Demokraten blev båret ind! Manden i lys kittel er avisforvalter A. Jacobsen. Demokraten 24.3.1918.*

# STATENS AVISSAMLING

- Oprettet 1916, åbnet 1918
- Samling opbygges i medfør af Pligtafleveringsloven
- Modtager 20.000 dagblade og 16.000 ugeaviser om året
- Samlingen omfatter i alt
  - 94.000 avisbind
  - 28.000 bind med distriktsblade
  - 66.000 mikrofilmspoler
  - 100 mio. sider i alt
  - 23.000 hylde-meter

# DET NYE MAGASIN FRA 2007



# DET NYE MAGASIN, INDEFRA



# DET NYE MAGASIN, INDEFRA





# DIGITALISERINGSPROJEKTET - BAGGRUND

# VESTINDISK PAKHUS

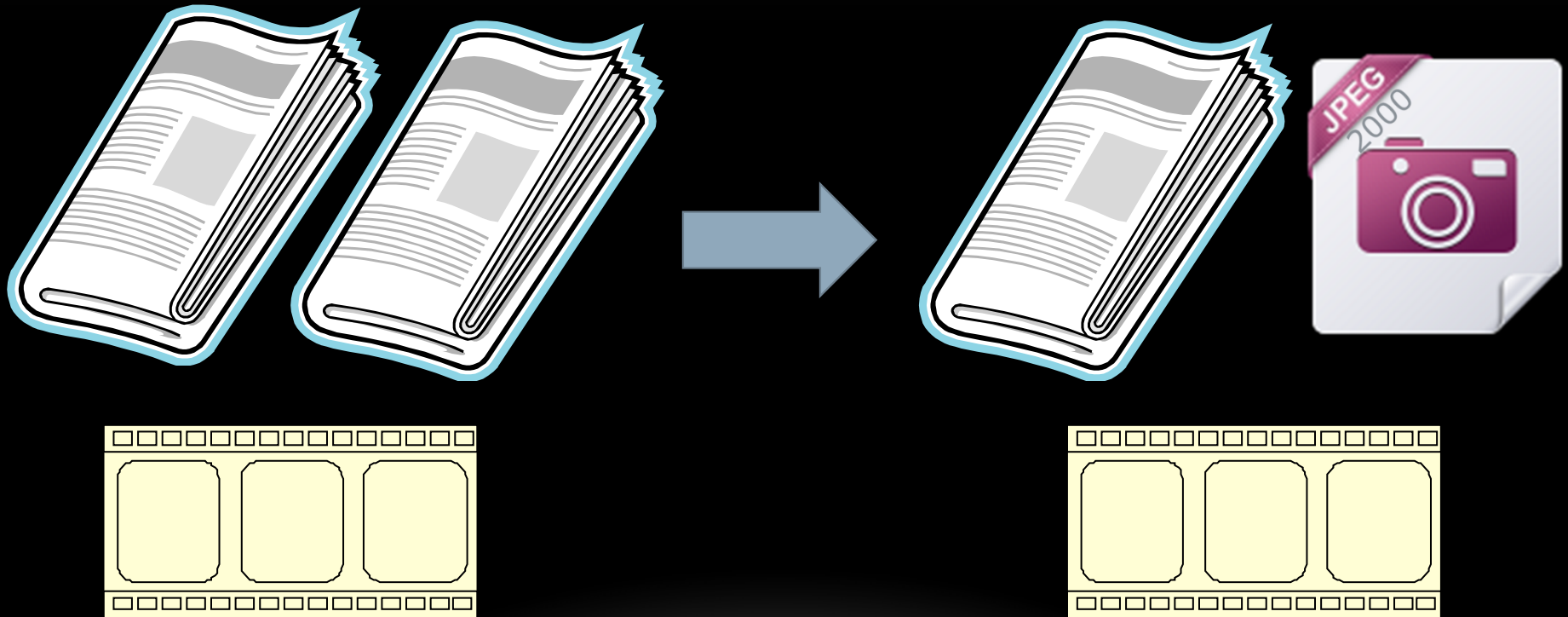
- KB har 32. mio. avissider liggende i Vestindisk Pakhus fra 1780, som ikke længere kan bruges til magasin.
- Avisdigitaliseringsprojektet finansieres ved at kassere disse aviser frem for at flytte dem til et nyt magasin



# FINANSIERING

- Bevaringsformål – Finanslov 2012 – 16 mio. kr.
  - Digitalisering med henblik på kassation
  - Genererer billedfiler
- Tilgængeliggørelse – UMTS midler – 7 mio. kr.
  - Brugergænseflade, søgning
  - OCR og segmentering
- Rigtig mange penge og arbejdstimer fra Statsbibliotekets løbende bevilling!

# STRATEGI FOR BEVARING

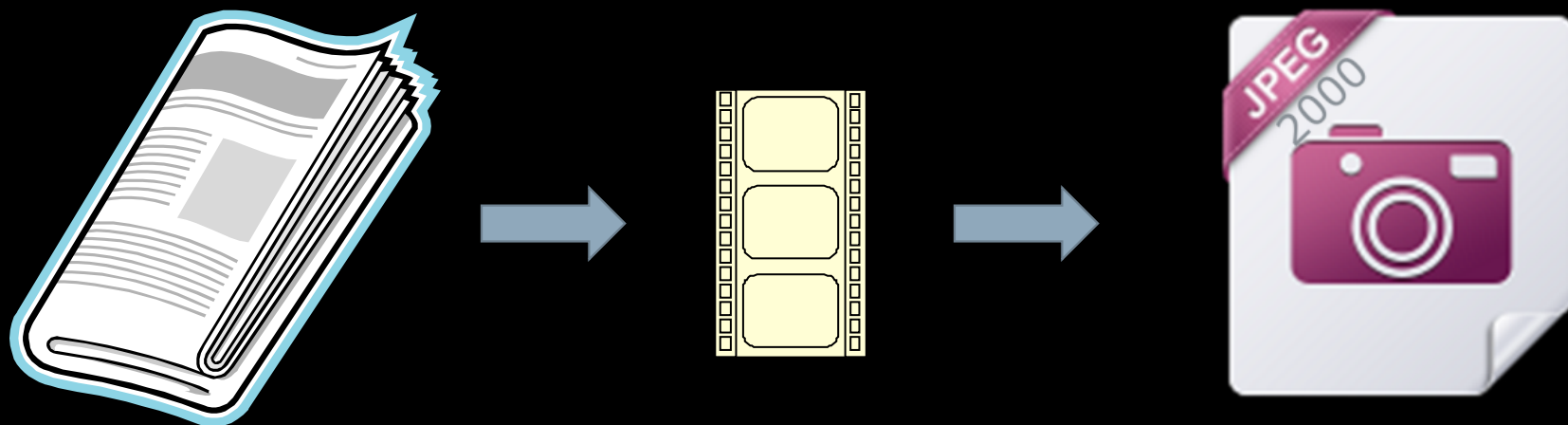


# TIDSPLAN FOR PROJEKTET

- 1. halvår 2014:
  - Pilotfase, produktionsmodning
- 2. halvår 2014:
  - Fuld produktion
  - 1.000.000 sider pr. måned
  - 50.000 sider pr. arbejdsdag
- 1. halvår 2017: Færdig med denne runde ...

# DIGITALISERING – HVORDAN?

# DIGITALISERING FRA MIKROFILM SOM VI ALLEREDE HAR



# MIKROFILMSPOLER





# I KASSER



# I REOLER



# TAL

- 1 kasse indeholder op til 15 æsker
- 1 æske indeholder 1 spole med 1 film
- 1 film indeholder i gennemsnit 500 billeder svarende til 1.000 sider
- Vi skal håndtere 32.000 spoler, godt 2.100 kasser, og alle de metadata der hører til
- Det bliver til 32.000.000 avissider med tilhørende metadata
- Op mod 140.000.000 filer
- 800 TB kulturarv

# HVILKE TITLER SKAL DIGITALISERES?

- Et stort puslespil, for at få mest mulig værdi for pengene
- De fleste nulevende dagblade
- Adresseavisen
- Aktuelt
- Land og folk
- Se listen på:  
<http://blog.avidigitalisering.dk/avistitler/>



# Ninestars

- Scanning i København hos Scanning.dk A/S
- Manuel efterbehandling i Chennai
- Hovedkvarter og IT-udvikling i Bangalore

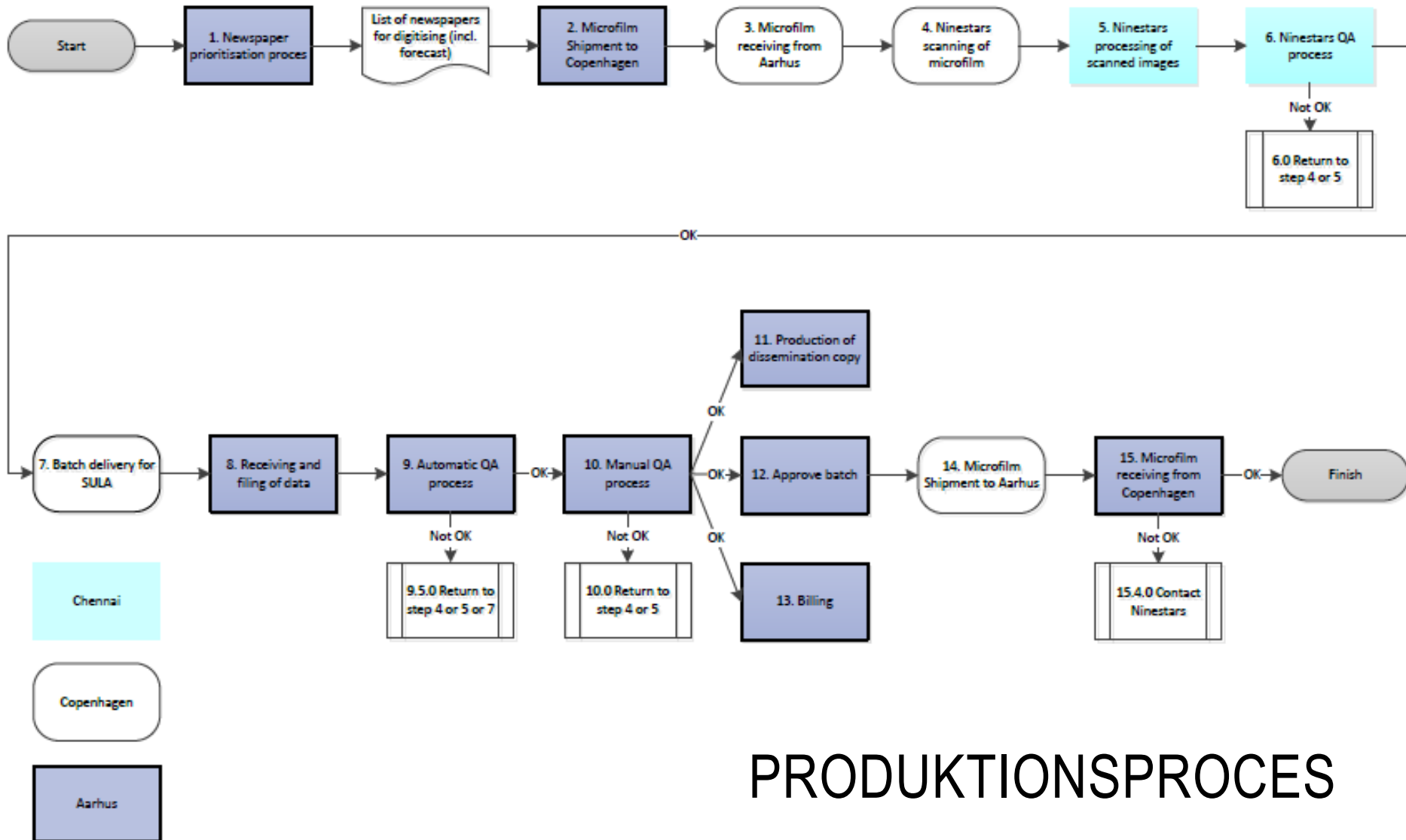


Ninestars bygning i Chennai



# Newspaper digitisation Top level production flow

Version 3.4  
Hamburg changed to Copenhagen



## PRODUKTIONSPROCES

# 1. PRIORITERINGSPROCES

- Hvad bestemmer rækkefølgen?
  - Kontrakter med de nulevende dagblade
  - Relevanskriterium

## 2/3. FILMENE SENDES TIL KØBENHAVN



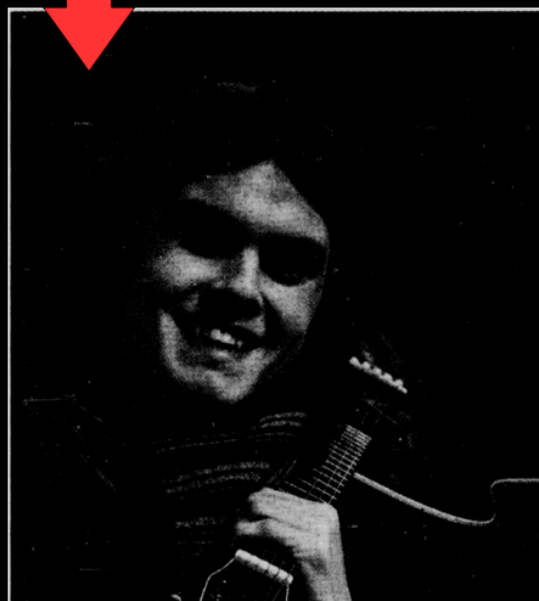
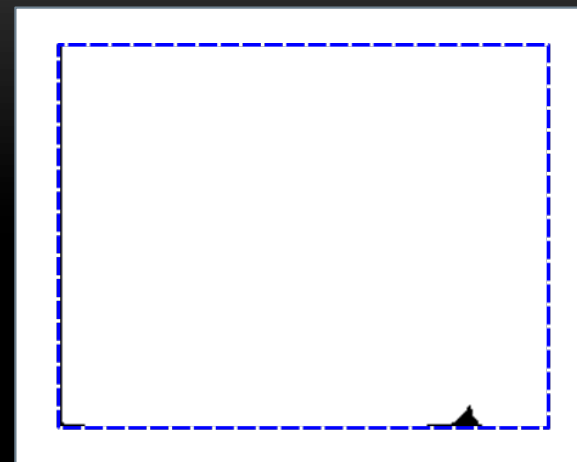


## 4. SCANNING AF MIKROFILM

- 2 Mekel Mach 5 scannere placeret i Albertslund
- Vi har været med ved kalibreringen af scannerne. Vi og Ninestars har ved fælles hjælp opnået et rigtig godt resultat.
- Det rå output er TIFF scannet i 300 dpi. Den endelige fil er en JPEG2000 og kæden mellem dem er lossless
- Vi får samtidig målt emulsion density på mikrofilmene til brug for kassationsprocessen



# BILLEDKVALITET – BEVARINGSKVALITET?



udtrykke sig.  
- På den måde kan jeg selv spille til, og er ikke afhængig af at skulle ud og finde en guitarist eller et band for at synge mine sange.

## Følsomme sange

Som man kunne se ved hans audition, så er det de mere følsomme toner, Mathias når guitareren findes frem.

- Jeg er mest til stille og rolig musik. Det er det, der tiltaler mig mest, men jeg har ikke nogen forbilleder som sådan siger han og tilføjer, at sangeren og sangskriveren Tim Christensen da er meget cool. Mathias skriver endnu ikke sange selv, men det er noget som han gerne vil kunne engang i fremtiden.



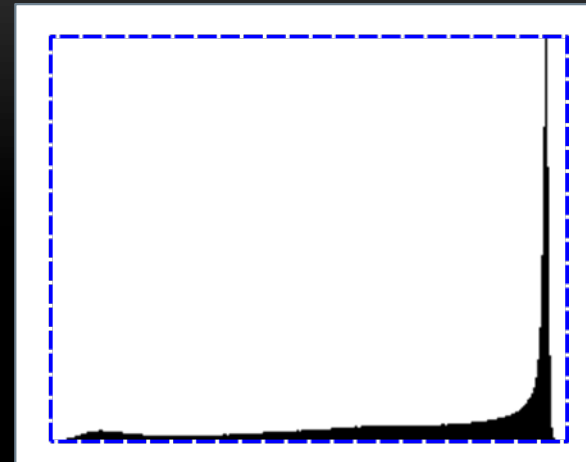
udtrykke sig.  
- På den måde kan jeg selv spille til, og er ikke afhængig af at skulle ud og finde en guitarist eller et band for at synge mine sange.

## Følsomme sange

Som man kunne se ved hans audition, så er det de mere følsomme toner, der tiltaler Mathias når guitareren findes frem.

- Jeg er mest til stille og rolig musik. Det er det, der tiltaler mig mest, men jeg har ikke nogen forbilleder som sådan siger han og tilføjer, at sangeren og sangskriveren Tim Christensen da er meget cool. Mathias skriver endnu ikke sange selv, men det er noget som han gerne vil kunne engang i fremtiden.

# BILLEDKVALITET – BEVARINGSKVALITET!



# 5. EFTERBEHANDLING I CHENNAI

- Split, crop og deskew
- Metadata
  - Tekniske metadata om den enkelte scannede side
    - Fx størrelse, opløsning, dato og operatør
  - Metadata om mikrofilmen
    - Titler, datoer, film og filmning
  - Udgave og udgivelse og den enkelte side
    - Morgen, aften, weekend, ekstra eller særudgave
    - Placering af siden i udgaven og på mikrofilmen
- OCR og segmentering



# OCR OG SEGMENTERING - FORMIDLING

(. 'Ilet I riipklIDii Jnryelt K||>  
(>! UliisN .il t..i\ci;<|rl l ' ! i :  
k 1 laagerup ni; f.irstni .1 m-  
'e 1 I | | i 'KP : ( Hl Jorgetl  
K]l»i»! I at l!!! I ! i i sll I '! I  
! : it I!' 1 dl; .len.-. Brum, i \  
ti . >l; ' r\ k .1 \ i iai nis  
Postens !•'< >nd J \ lut i n Is  
Posten A S , 'b'hn \ il>\ .1  
Telefon ffi kH bH kfi i- i n; s. i  
side

Chefredaktion: Jørgen Ejbøl  
(ansvarhavende), Ulrik  
Haagerup og Carsten Juste  
Direktion: Jørgen Ejbøl  
(administrerende) og Jens  
Bruun Udgiver og tryk:  
Jyllands-Postens Fond -  
Jyllands-Posten A/S 8260 Viby  
J. Telefon 87 38 38 38. Se  
også side 2

# 6/7/8/9/10 MODTAGELSE AF FILER OG KVALITETSKONTROL

- Ninestars sender de færdige batches til os, men inden da har de allerede kørt filerne igennem vores automatiske kvalitetskontrol.
- Vi modtager filerne, og kører dem igennem endnu en gang.
- Derefter foretager vi en manuel stikprøvekontrol og kigger på de sider, den automatiske kontrol fandt bemærkelsesværdige.

# MANUEL KVALITETSKONTROL

- 50.000 sider hver arbejdsdag
- ISO standard DS/ISO 2859
- Billedkvalitet – notere evt. fejkilde
- Kontrollere metadata
- OCR og segmentering





# HVILKEN KVALITET TIL HVILKET FORMÅL?

- Vi skal **bevare** de digitaliserede aviser i den bedst mulige kvalitet for fremtiden
- Vi skal **tilgængeliggøre** de digitaliserede aviser i den bedst mulige kvalitet for nutiden
- Vi skal kunne **afgøre** hvorvidt bevaringskopi og mikrofilm er tilstrækkelig gode til, at vi kan vælge at **kassere** den papirkopi, der ligger i Vestindisk Pakhus. For det er jo der, pengene kommer fra.

# NÅR DE DIGITALE FILER ER GODKENDT

- Fremstiller vi digitale formidlingskopier
- Modtager vi mikrofilmene retur fra København
- Begynder vi at arbejde med kassationsprocessen for en given titel
- Og vi tilføjer supplerende metadata

# SUPPLERENDE METADATA

- Sammenkædning af titler over tid
- Geografiske data
- Kobling til De Danske Aviser
  - <http://www.statsbiblioteket.dk/nationalbibliotek/adgang-til-samlingerne/aviser/de-danske-aviser>



avisID	titel	fra	til	DDA-id	udgivelsessted
avisID	title	from	to	DDA-id	place of publication
aktuelt	Socialisten	1871-07-21	1874-05-09	1-137	København
aktuelt	Social-Demokraten	1874-05-10	1959-03-31	1-137	København
aktuelt	Aktuelt	1959-04-01	1978-01-16	1-137	København
aktuelt	Lands-avisen Aktuelt	1978-01-17	1982-12-31	1-137	København
aktuelt	Aktuelt	1983-01-01	1987-05-02	1-137	København
aktuelt	Det Fri Aktuelt	1987-05-02	1997-08-07	1-137	København
aktuelt	Aktuelt	1997-08-08	2001-04-06	1-137	København

# TILGÆNGELIGGØRELSE

# TILGÆNGELIGGØRELSE

- Copyright, lovgivning, økonomi sætter rammerne.
- På Statsbiblioteket, Det Kongelige Bibliotek, Det Danske Filminstitut: adgang til alt via Mediestream/aviser.
- Online på Mediestream.dk: fri adgang for alle til aviser udenfor copyright.
- Adgang via nogle dagblades hjemmesider på varierende vilkår.
- Anden adgang måske senere?



# ADGANG TIL AVISERNE

På SB, KB, DFI

**Adgang til ALT**

Online for alle

**Adgang til gamle  
aviser**

Avisernes websites

**Adgang til  
enkelte avisers  
arkiv**

# ADGANG TIL AVISERNE

Periode	Adgang på SB og KB	Snippets/ thumbnails	Fjernadgang alle	Academic License
1950+	Ja	Afhænger af copydan-aftale		
1940	Ja			
1930	Ja			
1920	Ja			
1910	Ja			
1900	Ja	Ja	Ja	Ja
1890	Ja	Ja	Ja	Ja
1880	Ja	Ja	Ja	Ja
1870	Ja	Ja	Ja	Ja
1860	Ja	Ja	Ja	Ja
ældre	Ja	Ja	Ja	Ja



# SØGEMULIGHEDER

- Titel, dato og sidetal
- Tekstgenkendelse / OCR
  - Frakturskrift – lavere genkendelsesgrad
- Segmentering – identifikation af ”artikler” – automatisk

mediestream/aviser  
åbner foråret 2015

🏠 SØG I ALT RADIO TV REKLAMEFILM INFO LOG UD

# M=DI\_STREAM

Søg i Danmarks audiovisuelle kulturarv

Søg kun i det, du har adgang til.

### HVAD INDEHOLDER MEDIESTREAM?

Mediestream giver adgang til Statsbibliotekets digitale kulturarvssamlinger - indtil videre radio, tv og reklamefilm. Den næste store samling i Mediestream bliver 32 millioner avissider

### HVORDAN FÅR JEG ADGANG?

Alle kan søge i oplysninger om de enkelte udsendelser og reklamefilm. På Statsbiblioteket, Det Kongelige Bibliotek og Det Danske Filminstitut kan man se og lytte til materialet. Adgang fra andre steder kræver, at dit uddannelsessted har købt licens.

## DANMARKS AUDIOVISUELLE KULTURARV

f t

### Digitaliseret radio, tv og reklamefilm

Vi beriger hver dag radio- og tv-samlingerne med nye udsendelser, og samtidig arbejder vi på at udvide samlingen af reklamefilm med nye film fra tv og biograf. I 2014 åbner vi en ny samlingsindgang med ældre danske aviser.

**RADIO**  
FRA DANSKE KANALER

**REKLAMEFILM**  
FRA TV OG BIOGRAF

**TV**  
FRA DANSKE KANALER

**PÅ VEJ...**  
32 MIO. AVISSIDER

De

til Forsendelse med Posten

allene privilegerede

R i o b e n h a v n s k e

Maanedlige Stats=Lidender

Første Stykke

for

Januarii Maaned 1787.

Af disse Lidender udgives Maanedlig et Stykke, ved Brodrene Verling.

Den 22de Martii 1786 har Hs. Kongel. Majestæt aller-  
naadigst behaget at meddeele Amtmand Olav Stephens-  
sen i Island, for hans Godbædighed mod Trængende i de sam-  
mefeds sidst intrusne haarde Naringer, samt andre af ham  
tilforn ved adskillige Leiligheder udøvede patriotiske Handlin-  
ger, den store Medaille pro Meritis, i Guld, og tillige be-

Den 3die November er den deputerede Borger Christian  
Asmussen i Husum beskiftet til Raadmand; og Candidatus  
Theologiae Broder Brodersen kaldet til Pastor for Fabretoft  
i Amtet Tøndern, samt Candidatus Theologiae Kettel Bahns-  
sen til Pastor i Niesum i forberemeldre Amt. Under samme  
Dato er Doctor Medicinæ Fridrich Wilhelm Koch beskiftet

Den 3die November er den deputerede Borger Christian Asmussen i Husum bestikket til Raadmand; og Candidatus Theologiae Broder Brodersen kaldet til Pastor for Fahretoft i Amtet Tøndern, samt Candidatus Theologiae Kettel Bahnsen til Pastor i Riesum i forbemeldte Amt. Under samme Dato er Doctor Medicinæ Fridrich Wilhelm Koch bestikket til Physicus i Staderne Glückstad, Itzehoe, Crempe og Wilster, samt Amtet Steinburg.

Den 7de November ere Candidatus Juris Johann Siegra

Den 3die November er den deputerede Borger Christian Asmussen i Husum bestikket til Raadmand; og Candidatus Theologiae Broder Brodersen kaldet til Pastor for Fahretoft i Amtet Tøndern, samt Candidatus Theologiae Kettel Bahnsen til Pastor i Riesum i forbemeldte Amt. Under samme Dato er Doctor Medicinæ Fridrich Wilhelm Koch bestikket til Physicus i Staderne Glückstad, Itzehoe, Crempe og Wilster, samt Amtet Steinburg. Den 7de

Den 3die November er den deputerede Borger Christian Asmussen i Husum bestikket til Raadmand; og Candidatus Theologiae Broder Brodersen kaldet til Pastor for Fahretoft i Amtet Tøndern, samt Candidatus Theologiae Kettel Bahnsen til Pastor i Riesum i forbemeldte Amt. Under samme dato er Doctor Medicinæ Fridrich Wilhelm Koch bestikket til Physicus i Staderne Glückstad, Itzehoe, Crempe og Wilster, samt Amtet Steinburg.

# STATUS MEDIO NOVEMBER 2014

- Er digitaliseret:
    - Bornholms Tidende
    - Jyllands-Posten
    - Aarhus Stiftstidende
  - Er i proces:
    - Horsens Folkeblad
    - Adresseavisen
- i alt ca. 2.500.000 sider

FØLG MED OG LÆS MERE PÅ [avidigitalisering.dk](http://avidigitalisering.dk)

# HVAD SÅ NÅR VI ER FÆRDIGE MED DE 32.000.000 SIDER?

- Vi vil gerne lave crowdsourcing til forbedring af OCR, segmentering og metadata.
- Og der er stadig 30 mio. sider på mikrofilm og 40 mio. på papir at tage fat på.

SPØRGSMÅL?