

Testing Transkribus: Maskinassistert håndskriftgjenkjenning for små ABM-institusjoner

Frank Meyer, Næs Jernverksmuseum

Sammendrag

Denne rapporten beskriver testing av Transkribus, et verktøy for maskinassistert håndskriftgjenkjenning (HTR), på små samlinger av historiske dokumenter fra Næs Jernverksmuseum. Hovedmålet var å utforske hvordan ulike modeller kunne trenes på plattformen for å transkribere håndskrevne dokumenter automatisk.

Det var Næs Jernverksmuseum som tok initiativ til og søkte om støtte fra Arkivverkets utviklingsmidler for arkivsektoren i 2023. Forsker og historiker Frank Meyer var prosjektleder, Karl Richard Bie transskribent og Lars Tore Vintermyr prosjektmedarbeider med det operative ansvaret for den tekniske gjennomføringen. I løpet av sommeren 2024 tok prosjektleder kontakt med fire ressurspersoner fra Arkivverket, som selv hadde arbeidet med maskinassistert gjenkjenning av Jacob Aalls håndskrift, basert på arkivmateriale fra Arkivverket. Disse var Terje Brånå, André Nilsson Dannevig, Ingvill Marie Johansen og Maren Dahle Lauten. Til sammen drøftet prosjekt- og ressursgruppa framgangsmåte og resultater. Drøftingene var svært verdifulle og nokså omforente.

Arbeidet har gitt verdifull innsikt i hvordan Transkribus kan brukes for å automatisk gjenkjenne håndskrift i små dokumentsamlinger, men peker også på noen begrensninger i gratisversjonen og problemer med layout-gjenkjenning.

De viktigste hovedpunktene fra rapporten er:

1. Modelltrening i Transkribus:

- a. Det ble trent elleve modeller med Transkribus. De første fem modellene ble trent på et mindre datasett, mens de siste seks modellene ble trent etter hvert som mer manuelt transkribert materiale ble tilgjengelig. Noen modeller ble trent i kombinasjon med en eksisterende, større modell.
- b. Modeller ble trent ved å laste opp skannede sider, lage manuelle transkripsjoner og trene modeller på dette. Det var også nødvendig å rette opp feil i layouten, altså plasseringen til teksten, manuelt før modelltrening.
- c. Parametere som antall iterasjoner¹ og bruk av valideringssett² påvirket kvaliteten på de maskinskapte transkripsjonene.

2. Resultater fra første modellgenerasjon:

- a. Testene viste forbedringsmuligheter for modellene, spesielt ved bruk av flere iterasjoner og ved å ekskludere baksider av dokumenter.
- b. En modell som kombinerte egne manuelle transkripsjoner med en eksisterende, større modell, ga best resultater.

¹ En iterasjon er antallet ganger en modell har sett hele datasettet. En iterasjon kalles også gjerne en «epoch» eller «epoke».

² Et valideringssett er en del av dataene, spesifikt satt av til å teste en modell. Disse dataene er ikke modellen trent på.

3. Tekniske erfaringer og videre arbeid:

- Modeller med flere iterasjoner ga bedre resultater, men det ble også tydelig at optimalisering måtte skje innenfor et spesifikt antall iterasjoner (om lag 100–200).
- Det ble det testet videre med flere ulike blandingsmodeller³ i andre modellgenerasjon, og det var tydelig at en kombinasjon av egne manuelle transkripsjoner og en basemodell ga best resultater.
- Det var flere utfordringer, inkludert problemer med gjenkjenning av skrift på baksider av dokumenter og linjedeteksjon i enkelte tilfeller.
- Det er ønskelig å teste ut den beste modellen på flere tekster av Jacob Aall, for å få innsikt i hvor fleksibel modellen er.

4. Praktiske og økonomiske konklusjoner

- Det er overraskende hvor mye av Aalls håndskrift fra første halvdel av 1800-tallet Transkribus faktisk evner å transkribere.
- Når dette er sagt, er det likevel vanskelig å få noe mer enn et svært grovt inntrykk av hva innholdet på sida er, når man leser den maskinelle transkripsjonen. For mange viktige, meningsbærende ord (Prøver, Centrum, Herregård, Ovn mm) - for ikke å snakke om navn (Odense, Thejsen, Thomsen, Bentzen) - er ikke gjenkjent godt nok av Transkribus. Et eksempel er s. 409:

Maskinell transkripsjon	Manuell transkripsjon
Ved denne lighed tærker jeg at selde skrevige	Ved denne Leilighed tænker jeg at sende Hr. Thejsen
røveri af dskiblige, støæbte Sagen, som paa andre Stede	Prøver af adskillige støbte Sager, som paa andre Steder
finde god Afsætning. I Fommer har jeg begynds at et	finde god Afsætning. I Sommer har jeg begynds at faae
enbele hoste Kunder, som have afhedt en ikke ubetydetig	enkelte faste Kunder, som have afsadt en ikke ubetydelig
Deel af mit Støbegods, og jeg bilder mig ind at Odleste	Deel af mit Støbegods, og jeg bilder mig ind at Odense,
som or Caatum for, saa mange rege henegaarde og hos	som er Centrum for saa mange rige Herregaarde og Godse,
agsa kan afsætte m Del. Det er mig saaledes m	ogsaa kan afsætte en Deel. Det er mig saaledes om at
Gare at honrie ijige i Formngandel mas Migter, som	gjøre at komme igjen i Sammenhandel med Thejsen, som
masker Horsen vad mn alt hors vaske Frnrnfærs mi	som maaske Thomsen ved en alt for rask Fremfærd i min
Fraværelse stødte tilhede. Af Dig venter jeg at venstnet	Fraværelse stødte tilside. Af Dig venter jeg et venskablgt
Oed om Manden a at støe paa. Kommerser Bina	Ord om Manden er at stole paa. Kammerherre Bentzen
har varliger bestælt et korti Ane for sig og Andre,	har aarligen bestilt et Parti Ovne for sig og Andre,
og rimeligen maa altsaa Afsætning konne Give	og rimeligen maa altsaa afsætning kunne give
Nappe vil Bergs Gitter komme afstæd denne Gnig,	nogle til Bergi Gitter avstæd denne Gang,
esa jeg endnue ingen Tegning har faacti megtver	da jeg endnu ingen Tegning har faaet; men Ovnen

- En viss skuffelse ble forsterket av kjensgjerningen at museet i testopplegget hadde føret Transkribus med omtrent 70 manuelt transkriberte sider for å trene gjenkjenningsmodellen, og museet har brukt omtrent **kr 375 000,- i lønnsmidler** for testopplegget.
- Arkivverket testet håndskriftmodellen som prosjektgruppa ved Næs Jernverksmuseum hadde utviklet på eget materiale, og fikk likedan et skuffende resultat. Kolonnene nedenfor er henholdsvis sidenummer, modell, årstall og CER. Næs Jernverksmuseums «Jacob Aall – Handwriting» gjør det ikke bedre i testen enn de tidligere modellene som den er basert på. Dette er uventet, ettersom man ville tro at dere hadde lagd mer enn nok treningsdata for å få gode resultater på håndskriften til én enkel skribent. Det kan hende det bare trengs mer

³ En blandingsmodell er en modell som er trent på en blanding mellom egne treningsdata, altså transkripsjoner, og en eksisterende modell.

eksperimentering med modelltrening til, for å få et bedre resultat.⁴

5	19th century Danish Gothic handwriting	1815–1816	21,08	50,23
	18C Danish Administrative Writing (PyLaia)		25,32	56,81
	NorHand 1820–1940		21,08	58,69
	Jacob Aall - Handwriting		20,63	57,28
6	19th century Danish Gothic handwriting	1815–1816	14,61	46,43
	18C Danish Administrative Writing (PyLaia)		20,25	46,94
	NorHand 1820–1940		17,00	51,53
	Jacob Aall - Handwriting		14,71	44,39
7	19th century Danish Gothic handwriting	1815–1816	22,27	44,44
	NorHand 1820–1940		22,37	59,26
	Jacob Aall - Handwriting		27,21	65,28

- e. I det praktiske arbeidet kan den maskinelle transkripsjonen godt fungere som inspirasjon og korrektiv når leseren trenger en alternativ lesing av et ord eller en formulering som er vanskelig å tyde. Den kan på denne måten lette en manuell transkripsjon; men den er ikke i nærheten av å kunne erstatte denne.

⁴ E-post André Nilsson Dannevig – Frank Meyer, 30.10.2024.

Den detaljerte prosjektrapporten

Lars Tore Vintermyr, Frank Meyer

Manuell transkripsjon

Utvalget av Aalls brev kan begrunnes med at Aall var en utrøttelig skribent og en allsidig intellektuell, nærlingslivsaktør og politiker. Både i offentlige og private depotinstitusjoner fins det mange tusen sider med hans håndskrift. Mye av dette materialet kan sies å være viktig for å forstå tida Aall levde i, blant annet de dramatiske hendelsene som førte til Norges løsrivelse fra Danmark i 1814. På denne måten kan en transkripsjonsmodell være av stor nytte for nye generasjoner av historikere og andre interesserte.

Det ble alt i alt transkribert 165 sider av Jacob Aalls korrespondanse med sin beste venn Nils Hofmann Bang fra årene 1804 til 1810. De transkriberte sidene er en del av en større, sammenhengende korpus med originalbrev som Jacob Aalls etterkommere fikk i gave av etterkommerne til Nils Hofmann Bang. Til sammen utgjør korrespondansen omtrent 1300 sider og dekker tida fra 1804 til 1844, det vil si 40 år. Brevene er delt inn i to bunker som dekker tidsspennene 1804 til 1823 og 1824 til 1844 henholdsvis. Alle brev er digitalisert og tilgjengeliggjort på digitalarkivet:

- <https://media.digitalarkivet.no/view/150323/1>
- <https://media.digitalarkivet.no/view/150324/1>.

Aalls skrift er lett gjenkjennelig og skrevet i gotisk skript. Skriften er konsekvent og, med enkelte unntak, lett leselig.

Modell-trening i Transkribus

Prosjektmedarbeider Lars Tore Vintermyr har til sammen trent elleve («relativt distinkte») modeller, hvor de første fem modellene bygger på side 1–62 (63) av Næs Jernverksmuseums dokumenter⁵. De seks resterende modellene er bygget på samme samling, men etter at mer manuelt transkribert materiale ble tilgjengelig (s. 1–124). Antall modeller – at det er elleve og ikke flere – er en følge av begrensninger i gratisversjonen. Gratisversjonen av applikasjonen setter en grense på fem modeller som man kan trene i måneden.

Generell fremgangsmåte for modelltrening i Transkribus

Før en modell kan trenes i Transkribus-appen, må et utvalg skannede sider som inneholder håndskriften som modellen skal gjenkjenne, samt en manuelt transkribert versjon av disse sidene, lastes opp. Deretter velges en layout-modell som skal gjenkjenne hvor skriftlinjene er lokalisert i dokumentene, for så å nummerere hver linje den finner. Etter at layout-modellen har blitt brukt, vil linjene den har funnet som oftest inneholde noen mindre feil. Brukergrensesnittet i Transkribus gjør det enkelt å rette opp i disse feilene manuelt.

Brukergrensesnittet gir et intuitivt 1:1-forhold mellom den digitale faksimilen (bildefilen) og den maskinelle transkripsjonen (layout-modellen). Det skannede dokumentet ligger ved siden av et åpent dokument. Dokumentet inneholder bare nummereringen av linjene, slik at feilene kan bli

⁵ NESJ/NJM-005: Familien Aalls privatarkiv, E-00001: Korrespondanse Jacob Aall og Niels Hofman-Bang, L0001: Korrespondanse Jacob Aall og Niels Hofman-Bang, del 1c (<https://www.digitalarkivet.no/source/150323>).

rettet og de tilsvarende transkriberte linjene kan legges inn. Først når eventuelle feil og «støy»⁶ har blitt fjernet fra datasettet, er det mulig å trene en tekstgjenkjenningsmodell.

Når datasettet er ferdigstilt, kan man velge parametere som styrer treningsopplegget for den nye tekst-gjenkjenningsmodellen. Parameterne er (1) antall iterasjoner, (2) tidlig-stopp-funksjon, og i dette tilfellet vil også (3) valideringssettet ha en viss innvirkning på modellen, fordi datasettet er relativt lite.

- 1 Iterasjon-parameteren tilsvarer antall ganger algoritmen vil gå gjennom hele datasettet for å lære av det.
- 2 Tidlig-stopp-funksjonen stopper etter et gitt antall iterasjoner. Dette blir brukt for å unngå såkalt «overfitting». Det vil si at modellen stopper å ta til seg for mye informasjon som er spesifikt for datasettet, og som igjen vil begrense kvaliteten på transkribering av ny data.
- 3 Valideringssettet tilsvarer en andel av datasettet som blir brukt til validering av modellen. Denne delen av datasettet blir trukket ut før modellen blir trent. Algoritmen prøver å gjenskape mønstre den finner i datasettet, for deretter å sammenligne disse med valideringssettet. Den vil deretter omjustere («lære») hva som blir lagt vekt på, slik at den kommer så nær valideringssettet som mulig.
Ut fra denne sammenligningen vil treningsprosessen gi en CER (character error rate), som er en måleenhet for hvor mange tegnfeil den finner.

Til slutt vil alt dette bli lagret som en tekst-gjenkjenningsmodell. I tillegg kan man også bygge videre på etablerte modeller, og i tilfeller med små datasett kan dette være ideelt.

Konklusjoner etter den første modelltreningen

Den første av generasjonen av modeller hadde først og fremst som mål å finne området de ideelle parameterne lå i, slik at en mer raffinert versjon kunne bli trent i den neste generasjonen. Selv om testingen ikke ga noe definitivt svar, viste den en rekke forbedringsområder. Blant annet ble datasettet til den neste generasjonen av modeller mer begrenset, ved at de ikke inneholdt baksider.

Generelt presterte modellene som var basert på rundt 100 iterasjoner best, både med eller uten NorHand-modellen. Selv om de presterte bedre enn de andre modellene, **var** disse variasjonene designet til å være ekstreme for å gi en indikasjon på områdene modell-treningen burde utvikles i. Selv om de presterte best i eksempelet som er lagt til her, var det andre områder hvor de ikke presterte like godt. Ut fra dataene som ble samlet så langt, så det ut til at en modell med en minimal økning i iterasjoner, muligens basert på NorHand, kunne være ideell. Disse parameterne ble derfor prioritert i den neste generasjonen.

En annen fremgangsmåte prosjektgruppa ønsket å teste, var å trene opp en modell på deler av datasettet, for deretter å bygge videre på denne modellen ved å lage en ny modell. Denne nye modellen ville ha vært basert på den første modellen og de resterende sidene. I tillegg ville dette kunne gi en ny variabel ved å teste variasjoner i parameterne på de to modellene. Denne fremgangsmåten ble midlertidig satt til side, som følge av antallet modeller som måtte ha blitt lært opp. Hvis to modeller, en med og en uten NorHand skulle ha blitt trent, ville dette kreve minimum 4 unike modeller. (To basis-modeller og to påbygg-modeller, da tar man heller ikke hensyn til

⁶ Eventuelle feil gjort av layout-modellen samt annet som kan påvirke modelltreningen negativt, gjennomskinn, rifter i papiret med mer.

variasjoner i parameterne for hver modell.) Dette ville ha begrenset hvor mye annet prosjektgruppa kunne ha testet samtidig, som følge av 5-modell-begrensningen. Men med tilgang på et større antall manuelt transkriberte sider, testet prosjektgruppa denne fremgangsmåten i den andre generasjonen, og sammenlignet den med modeller trent på alt transkribert materiale.

Tabell 1: Første generasjon Jacob Aall-modeller rangert etter feilprosent (character error rate, CER).

Modell-navn	Epochs	Tidlig-stopp	Baksider	Bygget på Norhand	Inneholder siste side	Sammenligning av CER/WER	CER
Aall - med Norhand - lav epoch	60	20	ja	ja	nei	53.57 % / 90.05 %	48,80 %
Aall - handwriting - 1803-1823 NHand build v0.1	100	20	ja	ja	ja	20.02 % / 47.96 %	28,70 %
Model AallNoHand - v0.2	60	20	nei	ja	nei	28.27 % / 59.28 %	27.20 %
J.Aall - handwriting - v0.1	100	20	ja	nei	ja	20.95 % / 49.77 %	23.30 %
Aall - norhand - høy epoch (200) - med baksider	200	20	nei	ja	nei	17.52 % / 46.61 %	16.00 %

Lavere CER skal tilsvare færre feil i den maskinelle transkripsjonen, men (som nevnt på møtet med Arkivverket) er ikke denne indikatoren akkurat til å stole på.

Modell-trening av 1-generasjonsmodeller

Forskjellene i datasettet er begrenset til bruken av baksider av brevene (der de er inkludert), og en feil som ble gjort ved å legge inn siste side (63), som ikke var helt ferdig transkribert. Selv om det ikke er vist i eksempelet, så førte denne feilen til at de to modellene (overraskende og uventet) hoppet over deler av transkripsjonen på andre sider. Da problemet ble studert nærmere og modellene sammenlignet, var det tydelig at hyppigheten av dette problemet var mer markant i én av modellene, «J. Aall – handwriting – v0.1». Dette kom nok som følge av at denne modellen var mindre robust, ettersom den bare var basert på datasettet, og ikke bygget videre på NorHand, slik som Aall – handwriting – 1803-1823 Nhand build v0.1». Side 63 ville derfor ha tilsvart en større prosent-andel av datasettet i denne modellen og dermed hatt en større påvirkningskraft på treningen av modellen. Transkribus inneholder en funksjon som gir brukere muligheten til å «tagge» sider som «ikke ferdig transkribert», men i testopplegget valgte prosjektgruppa rett og slett å utelate den uferdige siden.

Problemet med baksider

En del av de digitale filene viste skrift på baksider. Layout-modellen («Universal Lines», også nevnt i lenken i neste avsnitt) kunne i disse tilfellene ikke gjenkjenne linjer som var skrevet sidelengs. Selv om baksidene ble rotert og layouten lagt inn manuelt, var ikke sidene rotert da layout-modellen gikk gjennom datasettet. Problemet ved løst ved en manuell redigeringen før treningen av tekst-gjenkjenningsmodellene startet.

Bilde 1: Skrift på baksiden av et brev



For å undersøke om dette hadde stor innvirkning på modellene, ble sidene med rotert eller sidelengs skrift fjernet i treningen av Model «AallNoHand – v0.2». Siden denne modellen i tillegg ble trent med færre iterasjoner er det vanskelig å vurdere hvor stor innvirkning den roterte skriften hadde på de maskinassisterte transkripsjonene. Men siden inkluderingen av baksidene i «Aall - med Norhand - lav epoch», som også er basert på 60 iterasjoner, presterer mye dårligere, er det antatt at denne variasjonen kommer fra baksidene. (Denne forbedringen gjelder også modellen med 200 iterasjoner som heller ikke inneholdt baksidene.) Ekskluderingen av disse sidene ble derfor gjort konsekvent i andre generasjons-modellene.

Hyperparametere i Transkribus

Modelltreningen i Transkribus-appen er noe begrenset av hvor mange parametere man kan endre på. I hovedsak er det tre parametere:

- 1 Epochs (iterasjoner – antall ganger applikasjonen gjennomgår treningsdataen)
- 2 En tidlig-stopp-funksjon som skal begrense «overfitting» på treningsdataene. (Det har ikke blitt observert noen mulighet til å begrense tidlig-stopp-funksjonen til mindre enn 20 steg.)
- 3 Prosentandelen av treningsdataene som ønskes å bli brukt som et valideringssett, hvor man kan velge mellom maskinelt valgte 2, 5, eller 10 prosent, eller «egendefinert», hvor man selv velger spesifikke sider som valideringssett. I testopplegget har 10 prosent blitt beholdt i alle tilfellene.

Siden antallet modeller er begrenset til fem i prøveversjonen, og fordi to modeller allerede hadde blitt lagd, ble det konsultert med devs og AI-entusiaster på forskjellige nettsider (stackoverflow/discord/quora/reddit/med flere) før treningen av de resterende modellene ble igangsatt. Selv om ingen av entusiastene hadde direkte kjennskap til Transkribus eller PyLaia-rammeverket som Transkribus er basert på, var det flere som hadde erfaring med liknende rammeverk som baserte seg på recurrent neural network (RNN). Prosjektmedarbeideren tok derfor spesielt hensyn til rådene som refererte til Kraken-rammeverket ettersom det er relativt likt til PyLaia, og var tidligere brukt i Transkribus samt eScriptorium (et annet verktøy for maskinassistert håndskriftsgjenkjenning).

Et av punktene som ble nevnt på brukerforumet var at 100 iterasjoner var langt mer enn hva et datasett på rundt 6000 ord skulle tilsi (antallet ord i datasettet som ble brukt her). Det var derfor sannsynlig at problemene med modellene kom som følge av «overfitting», og at en test på langt færre iterasjoner burde bli gjort. Ettersom det ikke er mulig å begrense antallet i stopp-funksjonen og fordi rundt seks sider ble brukt til valideringssettet, ble antallet iterasjoner begrenset til 60 i modellene «Model AllNoHand – v0.2» og «Aall - med Norhand - lav epoch».

(Sammelijning av kraken og PyLaia: HTR+ versus PyLaia (uni-greifswald.de))

Det ble derimot fort tydelig at disse to modellene presterte langt dårligere enn de to første, og som en ren test, ble det derfor trent en modell med dobbelt så mange gjennomganger som originalene (200). Den presterte bedre på test-siden enn 60-modellene, men som følge av at 100-modellene presterte relativt like godt på den inkluderte test-siden, som begge inneholder et dårligere datasett, indikerte det at en ideell modell sannsynligvis ville ligge i intervallet 100–200 iterasjoner, og trolig ikke over 200. Som nevnt er det ikke nødvendigvis slik at man kan stole på det innebygde CER-verktøyet. Ut fra den enkeltsiden som er lagt til her, ligger nok den beste transkripsjonen et sted mellom «J.Aall - handwriting – v0.1» og «Aall - handwriting – 1803–1823 NHand build v0.1». Det viste seg at «J.Aall - handwriting – v0.1» var best på særtrekk ved teksten, mens «Aall - handwriting – 1803–1823 NHand build v0.1» var best på hverdagspråk som følge av NorHand. Som nevnt, ble en versjon med litt høyere antall iterasjoner, uten siste side, uten baksidene og som baserte seg på NorHand, trent i den andre generasjonen med modeller som var nærmere en slik likevekt.

Modell-trening av 2-generasjonsmodeller

Denne testen var i første omgang begrenset til en sammenlijning av to modeller med identiske parametere (iterasjoner, blandingsmodell og identisk datasett). Det som skiller dem, er hvordan datasettet har blitt brukt i treningen. Modellen «JAall + Norhand stor v.02» er trent på vanlig måte og er bygget på alt manuelt transkribert materiale i en enkelt treningsprosess.

I treningen av modellen «JAall + Norhand split v0.1» ble derimot datasettet delt i to. En modell ble trent på den første delen av datasettet, splitt-modellen er basert på denne modellen og den resterende delen av datasettet. I tillegg inneholdt testen en outlier kalt «JAall - 150 – v0.1». Denne tilsvarte «J.Aall – handwriting – v0.1» fra forrige test, med en liten økning i antall iterasjoner og et større datasett. Denne modellen skulle først og fremst se på effekten av en blandingsmodell når datasettet var dobbelt så stort.

I møtet med Arkivverket ble det gitt mange gode råd om muligheter til å forbedre modellene på. Blant annet ble det vist til et effektivt sammenlijningsverktøy som var tilgjengelig i desktop-versjonen av Transkribus, men som ikke var tilgjengelig i appen. Det ble også henvist til tre modeller som hadde gitt gode resultater og som kunne bli brukt som blandingsmodeller.

De tre modellene det ble henvist til var:

1. «18C Danish Administrative Writing (PyLaia)»
2. «19th century Danish Gothic handwriting»
3. «NorHand 1820–1940».

Som følge av begrensingen til fem modeller, ble disse tre modellene testet på valideringssiden (s. 71) (vedlegg) og sammenlijnet ved hjelp av sammenlijningsverktøyet. I det tilfellet denne sammenlijningen ville vise noen overlegen modell, selv om II og III ser ut til å ha prestert best, var

det ikke noen tydelig forskjell i modellene. Siden NorHand ble brukt i modellene som ble trent i forrige test, ble Norhand-modellen prioritert for en mulig sammenligning med disse, men en test av de to andre modellene var fortsatt ønskelig. «19th century Danish Gothic handwriting» ble trukket frem som en spesielt god modell, og en sammenligning av modellenes parametere viste at antallet iterasjoner brukt i «19th century Danish Gothic handwriting» lå nærmere de ideelle verdiene prosjektgruppa hadde kommet frem til. Generelt presterer de tre nye Norhand-baserte modellene bedre enn de forrige fem, men effektiviteten er fortsatt veldig begrenset. Ettersom denne testen ble gjort i månedsskiftet og uten noen variasjon i datasettet, ble også modeller basert på de to andre blandingsmodellene testet. I tillegg ble en modell med en liten økning i iterasjoner basert på «19th century Danish Gothic handwriting» trent som følge av at «JAall + Danish Handwriting(130) v0.1» presterte bedre enn de andre modellene.

Tabell 1: modeller rangert etter størst feilprosent

Modell-navn	Iterasjoner	Sammenligning av CER/WER %	Bygget på modell:	CER %
JAall - 150 - v0.1	150	29.95/65.16	ingen	19.00
JAall + Norhand split v0.1	130	44.53/79.35	NorHand	15.00
JAall + Norhand stor v.02	130	23.05/ 48.39	NorHand	10.00
JAall + 18C Danish Administrative Writing v0.1	130	20.44/ 42.58	18C Danish Administrative Writing	10.00
JAall + Danish Handwriting (130) v0.1	130	17.32/ 34.84	19th century Danish Gothic handwriting	9.00
JAall + Danish Handwriting (150) v0.1	150	17.97/ 37.42	19th century Danish Gothic handwriting	7.00

**tidlig-stopp er fjernet.*

Konklusjon

Som følge av de to testene er det nå langt færre områder som kan optimaliseres, og det er i hovedsak begrenset til antall iterasjoner og bruk av blandingsmodell⁷.

Basert på sammenligningen av de tre NorHand-modellenes transkripsjon av side 71, ser det ut til at «JAall + Norhand stor v.02» er den beste av disse tre, men den presterer fortsatt dårligere enn modellene basert på de andre blandingsmodellene. Det er tydelig at «JAall + Danish Handwriting(130) v0.1» og den tilsvarende 150-modellen kom best ut i testen, der hvor de andre modellene har lengre linjer med feil, er disse to i større grad begrenset til enkeltord. 130-modellen presterer minimalt bedre enn 150-modellen, det er to mulige forklaringer på dette. Den ene forklaringen er at det kommer som følge av valideringssettet, som nevnt blir 10 prosent av datasettet tilfeldig trukket ut, og 150-modellen kan derfor ha blitt gitt et dårligere valideringssett. For en mer nøyaktig sammenligning kunne et identisk valideringssett ha blitt trukket ut manuelt, men 150-modellen ble først og fremst trent for å se etter mulig optimalisering. Det var heller ingenting som indikerte at valideringssettet til 130-modellen var et spesielt godt valideringssett, og hvis 150-modellen hadde blitt trent for en sammenligning, ville dette ha utelukket variasjoner med bedre valideringssett. Den andre forklaringen som trolig er den mest sannsynlige, er at global

⁷ Se fotnote 3.

optima⁸ for modell-treningen ligger i intervallet 120-150 iterasjoner. Argumentet for dette er basert på hvor konsekvent modellene som er trent med 120-150 iterasjoner presterer best, både når det gjelder blandingsmodellene og modellene de prosjektgruppa har trent. Hvis dette er tilfellet, vil det kunne forklare hvorfor 130-modellen er mer nøyaktig, og at denne nøyaktigheten begynner å avta i 150-modellen fordi modellen ligger lengre unna global optima. Dette vil også forklare hvorfor den NorHand-baserte modellen gjorde det dårligere enn de andre to. Basert på hvordan modellen med 200 iterasjoner i den første generasjonen viste tegn til overfitting⁹, og at modellene med 60 iterasjoner presterte så dårlig, kan man utelukke at dette er et «local optima»¹⁰.

Splitt-modellen ville kanskje ha gitt bedre resultater med et større datasett, eventuelt høyere antall iterasjoner eller en variasjon i iterasjoner mellom de to delene. Dette kunne være veldig interessant å se på, men det ser ut til at den krever andre variabler enn de som har blitt brukt her, samtidig krever den også to av de fem mulige modellene som kan produseres, så det begrenser hvor mange splitt-modeller som kan produseres samtidig. Som følge av hvor dårlig modellen presterte, ble denne typen modeller lagt til side, og det ble lagt større vekt på modeller basert på de andre blandingsmodellene.

Selv om «18C Danish Administrative Writing (PyLaia)», kom dårligst ut i bruken av sammenligningsverktøyet, var det ikke nødvendigvis slik at en modell basert på denne ikke kunne konkurrere med de andre to. De tre blandingsmodellene som ble testet her, varierer både i størrelse og i antall iterasjoner som har blitt brukt i treningen. Datasettene ligger i intervallet 600 000 ord for den minste og 1,5 millioner for den største, og antallet iterasjoner varierer mellom 124 og 207. Til sammenligning ligger modellene som er trent her på litt under 12 000 ord. De mest generelle modellene eller modellene med liknende skrift, (eventuelt samme forfatter) i treningssettet kommer derfor best ut, men denne sammenligningen sa ikke noe om vektene til de forskjellige modellene, og det var derfor *teoretisk* mulig at en modell som presterte dårlig uten datasettet, kunne prestere bedre etter at den ble trent på datasettet. Dette viste seg tydelig i modellen som ble trent, selv om den presterte dårligere enn «19th century Danish Gothic handwriting»-modellene.

Problemer ved bruk av Transkribus-appen under testingen

Til slutt vil vi omtale tre problemer prosjektgruppa støtte på under testingen.

Problemer ved bruk av Transkribus-appen og modelltreningen (første generasjon)

Et lite tillegg som er verdt å nevne, er at det har vært noen problemer med bruk av appen. Den har ulike ganger vært treg til å laste inn ved å låse seg i innlasting (selv etter at man har gått ut av nettleseren, slettet temp-folder/cookies/historikk med mer, og starter opp igjen på nytt). Den har også en rekke ganger begrenset opplastingen av flere filer, hvor det har ført til en feilmelding, og filene måtte lastes opp igjen en og en. Det er uklart hvorfor dette skjedde. Men siden problemet gjentok seg både i Firefox og i Microsoft Edge (begge uten nettleser-tillegg), og testen fant sted i tidsrommet kl. 11.00–14.00, kan problemet muligens skyldes høy trafikk eller liknende fra Transkribus sin side. Det er også mulig at problemet utelukkende gjelder gratisversjonen.

⁸ Globalt optimum er det beste resultatet som kan oppnås.

⁹ «Overfitting» skjer når en modell blir overtrent på en del av dataene, og dermed blir dårligere på andre deler av dataene.

¹⁰ Lokalt optimum er det beste resultatet som kan oppnås innenfor en lokal begrensning.

Problemer ved bruk av Transkribus-appen og modelltreningen (andre generasjon)

Det var en bug som skapte litt problemer. Denne feilen oppstod sannsynligvis etter at en side ble fjernet fra datasettet. Dette førte til feilmeldingen «Cannot invoke "String.length...s null (NullPointerException)», hvor applikasjonen refererte til siden som ble slettet, men hvor referansen til siden fortsatt eksisterte i samlingen, så siden som ble slettet, forsvant ikke. Denne feilen ble så videreført til neste side og i tillegg påvirket den hele samlingen siden hadde vært del av. Denne feilen kan ha påvirket siste del av splitt-modellen, selv om splitt-modellen ikke inneholdt den aktuelle siden. Det er usikkert når denne feilen først oppstod, for den ble først oppdaget da samlingen ble kopiert. Feilen førte til at denne delen av datasettet måtte lastes opp og bearbejdes på nytt. Siden datasettet i denne testen var relativt lite, så var ikke dette et stort problem. Men som følge av feilen og den nye opplastingen, lurte side 71 seg inn i første iterasjon av «JAall + Norhand stor». Siden ble fjernet i v0.2, og ideelt ville en modell basert på «19th century Danish Gothic handwriting» blitt trent i stedet for denne, men som nevnt ble denne lagt til senere sammen med 18C Danish Administrative Writing (PyLaia).

Problem med gjenkjenning av linjer

Et problem som ikke påvirket testen, men som sannsynligvis vil være mer utbredt i videre bruk av Transkribus, er gjenkjenningen av linjer gjort av «Universal lines»-modellen¹¹. Den hadde en tendens til å finne linjer i midten av ordene. Dette har ikke påvirket modellene, fordi det ble rettet opp manuelt i datasettet der det oppstod. Men det vil muligens begrense effekten av applikasjonen betraktelig hvis dette må gjøres manuelt for hver side under automatisk transkribering, spesielt med hensyn til lønnsomheten ved bruk av KI som følge av tid brukt til manipulasjon av «miljøet» før det kan gi gode resultater.

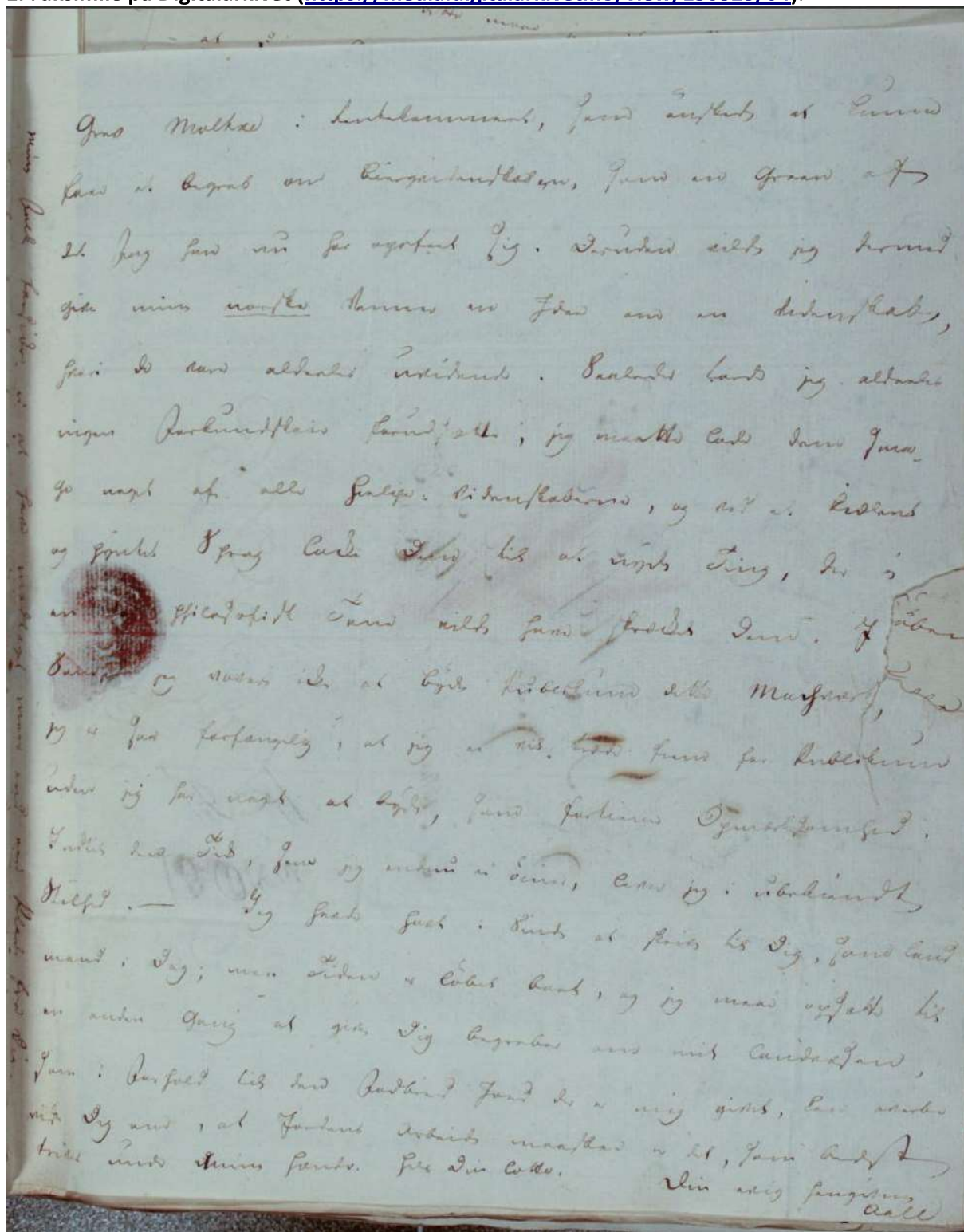
I møtet med Arkivverket og i en medfølgende rapport prosjektgruppa fikk tilgang på, ble det også lagt vekt på problemer med «Universal lines»-modellen. I Arkivverkets tilfelle var problemene begrenset til linjegjenkjenning av dobbeltsider, hvor linjene ble sammenhengende på tvers av dobbeltsidene. Ettersom materialet prosjektgruppa testet var begrenset til brev, har prosjektgruppa ikke støtt på dette problemet ennå. Men siden de ikke hadde det samme problemet som prosjektgruppa, er det mulige problemet vårt er begrenset til akkurat denne håndskriften.

Et annet aspekt ved denne testen, og grunnen til at side 71 ble valgt som valideringsside, var å undersøke i hvor stor grad avkuttete ord ville påvirke modellene (for eksempel «Høst v». «nedtake Ko». «til Kak»). Avskårne ord ble fjernet fra datasettet modellene er trent på, i frykt for at de ville påvirke modellene negativt, og som følge av det var testen litt «urettferdig». Derfor ble disse ordene ikke tillagt så mye vekt som resten av testen. Men ettersom det potensielt kunne vise hvor «rigide» modellene er (som følge av «overfitting»), var det noe vi ønsket å teste. Testen viste at modellene ikke hadde noen akutt forskjell, men kan ha hatt større interesse i sammenligningen av de tre blandingsmodellene, hvor det så ut til at «19th century Danish Gothic handwriting» kom best ut.

¹¹ «Universal Lines» er standardmodellen for layout i Transkribus.

Vedlegg: Manuelle og maskinassisterte transkripsjoner

1. Faksimile på Digitalarkivet (<https://media.digitalarkivet.no/view/150323/64>):



2. Manuell transkripsjon

- 1-1 # Grev Moltke i Rentekammeret, som ønsket at kunne
1-2 # faae et begreb om Biergvidenskab ere, som en Green af
1-3 # det Fag han nu har opofret(?) sig. Desuden vilde jeg dermed
1-4 # give mine norske Venner en Idee om en Videnskab,
1-5 # hvori de vare aldeeles uvidende. Saaledes torde jeg aldeeles
1-6 # ingen Forkundskaber forudsætter; jeg maatte lade dem sma-
1-7 # ge noget af alle Herlige-Videnskaberne, og ved et kiælent
1-8 # og pyntet Sprog lokker Den til at nyde Ting, der i
1-9 # en saa(?) filosofisk Tone vilde have skrækket Dem. I [*rift i siden]
1-10 # Sandheed(?) jeg vover ikke at byde Publikum dette Machværk,
1-11 # jeg er saa forfængelig, at jeg ei vil træde frem for Publikum
1-12 # uden jeg har noget at byde, som fortienet Opmærksomhed.
1-13 # Indtil den Tid, som jeg endnu ei øiner, lever jeg i Ubekiendt
1-14 # Stilhed. – Jeg havde havt i Sinde at skrive til Dig, som land-
1-15 # mand i Dag; men Tiden er løbet bort, og jeg maae opsætte til
1-16 # en anden Gang at give Dig begreber om mit landvæsen,
1-17 # som i Forhold til den Fodbred Jord der er mig givet, kan overbe[-]
1-18 # vise Dig om, at Jordens Arbeide maaskee er det, som Bedst
1-19 # trives under Diine(?) Hænder. Hils Din Lotte.
1-20 # Din evig hengivne
1-21 # Aall

3. Maskinassisterte transkripsjoner (Gen-1)

Her sammenlignes hver modells transkripsjon av side 63 med den manuelle transkripsjonen. Ord som er farget rødt og kryssset ut, er ord modellene har transkribert feil. Ord som er farget grønt er de tilsvarende ordene fra den manuelle transkripsjonen. Ord som er farget hvitt er ord som modellene har transkribert riktig.

3.1 Aall - norhand - høy epoch (200) - med bakside:

- 1-1 # ~~ger~~ Mavta Venkekommert, Grev Moltke i Rentekammeret, som ~~unkede a Komne~~ ønsket at kunne
1-2 # ~~faa~~ at begret on Bergendskaden, faae et begreb om Biergvidenskab ere, som en ~~Graan~~ Green af
1-3 # det ~~seg~~ Fag han nu har ~~opatret~~ opofret sig. ~~Deruden~~ Desuden vilde jeg dermed
1-4 # give mine ~~norkte Vmner~~ norske Venner en ~~hdre~~ Idee om en Videnskab
1-5 # ~~haa~~ hvori de vare aldeeles uvidende. ~~Saaleder~~ Saaledes torde jeg ~~aldele~~ aldeeles
1-6 # ingen ~~Forkundskale~~ forudsætte, Forkundskaber forudsætter; jeg maatte ~~lode~~ lade dem ~~sia~~ sma
1-7 # ~~ga~~ ge noget af ~~elle~~ ~~hjelge~~ Herlige Videnskaberne, alle Herlige Videnskaberne, og ved et ~~kiædens~~ kiælent
1-8 # og ~~hyntets~~ pyntet Sprog ~~lakke~~ Dilig lokker Den til at nyde Ting, ~~dereg~~ der i
1-9 # ~~Milæsafiste~~ filosofisk Tone ~~eilde~~ ~~haves~~ ~~srikket~~ ~~din~~ vilde have skrækket Dem. I
1-10 # ~~Saudyer~~ oeg ~~væver~~ Sandheed jeg vover ikke at byde ~~tublekomm~~ dete Madhvar. Publikum dette Machværk,
1-11 # ~~D~~ g jeg er saa ~~forfangelig~~, forfængelig, at jeg ~~a~~ vid Brode ~~fem~~ ei vil træde frem for ~~kublokun~~ Publikum
1-12 # ~~ader~~ uden jeg har noget at ~~bydi~~ byde, som ~~fortiene~~ Opmærksomhed, ~~fortienet~~ Opmærksomhed.
1-13 # ~~ende~~ Indtil den ~~rud~~ Tid, som jeg endnu ei ~~Kun~~ laver øiner, lever jeg i ~~ikekindt~~ Ubekiendt
1-14 # ~~Silhed~~ ne jeg Stilhed. – Jeg havde ~~haat~~ havt i Sinde at ~~frive~~ skrive til Dig, som ~~kande~~ land
1-15 # ~~maed~~ Dog mand i Dag; men Tiden er løbet bort, og jeg maae ~~opsalte~~ opsætte til
1-16 # ~~sen~~ en anden Gang at give Dig begreber om mit ~~Landvæsen~~ landvæsen,
1-17 # som i Forhold til den Fodbred ~~hard~~ de Jord der er mig givet, kan ~~øverke~~ overbe
1-18 # ~~gvise~~ vise Dig om, at ~~Fordens~~ Vrbeide ~~maarken~~ Jordens Arbeide maaskee er det, ~~hom~~ bedste som Bedst
1-19 # ~~srives~~ unde ~~Dine~~ hænder ~~hils~~ trives under Diine Hænder. Hils Din ~~Lote~~ Lotte.
1-20 # Din evig hengivne
1-21 # ~~h~~ Aall

3.2 J.Aall - handwriting - v0.1:

- 1-1 # Grev Mala Grev Moltke i tintelmert, Rentekammeret, som unsket-ønsket at kunne kunne
1-2 # faa at begrel faae et begreb om bergvidenskaln, Biergvidenskab ere, som en Green a-Green af
1-3 # det sog Fag han nu har opofferet sig i Druden vilde opofret sig. Desuden vilde jeg demed-dermed
1-4 # give min norste mine norske Venner en hde Idee om en Videnskat Videnskab
1-5 # paa hvori de var aldeales nvidens. Saalede vare aldeeles uvidende. Saaledes torde jeg alteles-aldeeles
1-6 # nigen Forkundskavr forindsatte, ingen Forkundskaber forudsætter; jeg maatte lade dem sa sma
1-7 # ga ge noget af alle hilge Viidenskabrie, Herlige Videnskaberne, og ved et kiellens kiælent
1-8 # og Gntet pyntet Sprog lakte Ditg lokker Den til at nyt Dig, de og nyde Ting, der i
1-9 # Ghlosofiste Tane alde hav hræktet dem, filosofisk Tone vilde have skrækket Dem. I
1-10 # Sadtter Sandheed jeg vver vover ikke at byde fublkine Publikum dette Mahvæs Machværk,
1-11 # jeg er saa forfangelg, forfængelig, at jeg a ei vil bæde fom træde frem for Vblatim Publikum
1-12 # ader uden jeg har noget noget at bydi, byde, som fortienne Opmækfomhed, fortiener Opmærksomhed,
1-13 # saller Indtil den gdis, Tid, som jeg endan er Tune, laver endnu ei oiner, lever jeg i ibekindt Ubekiendt
1-14 # Stilhed er jeg Stilhed. - Jeg havde haat havt i kinde Sinde at friv skrive til Dig, sam kan som land
1-15 # maad mand i Dag; men Tiden er kabot løbet bort, og jeg maae opsalte opsætte til
1-16 # en anden Gaig Gang at give Dig begrber begreber om mit landavsæn, mit landvæsen,
1-17 # som i Forhold til den Fdbred fored de Fodbred Jord der er mig givet, lar ovirb kan overbe
1-18 # giise vise Dig om, at fordens Vrvend Jordens Arbeide maaskee er det, som be be-Bedst
1-19 # ve trives under Vne hønder his Diine Hænder. Hils Din toke-Lotte.
1-20 # Din avig hangilig evig hengivne
1-21 # Aall

3.3 Aall - handwriting - 1803-1823 NHand build v0.1:

- 1-1 # Gea Maltas Grev Moltke i tentekmmert, Rentekammeret, som unsked-ønsket at kunne
1-2 # faa faae et begret begreb om keergvidenskalen, Biergvidenskab ere, som en Green as af
1-3 # det seg Fag han u nu har ogofret opofret sig. Dernden-Desuden vilde jeg dened-dermed
1-4 # give mine norste Venner norske Venner en hide Idee om en Videnskab, Videnskab
1-5 # havi hvori de vare aldeeles uvidende. Saalede Saaledes torde jeg aldreliis aldeeles
1-6 # nigen Forkundskave forindsatte, ingen Forkundskaber forudsætter; jeg maatte lade den soma dem sma
1-7 # gi ge noget af alle Hielge Vi denskabrmn, Herlige Videnskaberne, og ved et kiælens kiælent
1-8 # og hyntiet Sprag lakke Dilig pyntet Sprog lokker Den til at nyde Dig, de og Ting, der i
1-9 # ghilasafist Dane vilde havi Hrækket Den. Ie filosofisk Tone vilde have skrækket Dem. I
1-10 # Sausgger Sandheed jeg Næser vover ikke aet at byde Fublekine Publikum dette Mahver Machværk,
1-11 # jeg er saa forfangelg, forfængelig, at jeg et mud bræde som ei vil træde frem for Sublifinn Publikum
1-12 # nden uden jeg har, noget har noget at byidig, byde, som fortine Opmaelsomhed, fortiener Opmærksomhed,
1-13 # sndler Indtil den orkis, Tid, som jeg enden endnu ei Kun, laver oiner, lever jeg i iubekindt Ubekiendt
1-14 # Snlhed rer Stilhed. - Jeg havde haat havt i kinde Sinde at skrive til Dig, som lan land
1-15 # maed Deg, mand i Dag; men Tiden er kabot løbet bort, og jeg mee opsatte maae opsætte til
1-16 # en anden Ging Gang at give Dig begrelre momm met landvgæn, begreber om mit landvæsen,
1-17 # som i Frhold Forhold til den Fodbred saree Jord der er mig givet, kar avrker kan overbe
1-18 # givise vise Dig om, at Hordens ilibeide maaskeen Jordens Arbeide maaskee er det, som beids be Bedst
1-19 # res mende Dinne hande his trives under Diine Hænder. Hils Din Lole, Lotte.
1-20 # Din evig hengivn hengivne
1-21 # Ja Aall

3.4 Model AallNoHand - v0.2:

- 1-1 # Ga Maka Grev Moltke i Rintekammes, Rentekammeret, som ønsket at kunne kunne
1-2 # faa at begel m borgandenskalen, faae et begreb om Biergvidenskab ere, som en Green-Green af
1-3 # det seg Fag han nu nu har øgotret sig Deruden opofret sig. Desuden vilde jeg demed dermed
1-4 # give min morste Komer e heden o e Videnskatg mine norske Venner en Idee om en Videnskab
1-5 # haa hvori de vare aldeales alidende, aldeeles uvidende. Saaledes torde jeg aldendes aldeeles
1-6 # nigen Forkudskale forudsalte, ingen Forkundskaber forudsætter, jeg maatte lade den som dem sma
1-7 # gi nget ge noget af elle halge Vdenskabr, alle Herlige Videnskaberne, og ved es killens et kiælent
1-8 # og givnetet Sjrig likke Dilig pyntet Sprog lokker Den til at nyede Dig, de og nyde Ting, der i
1-9 # Elasaalste Due alde hane i hrllket Den, filosofisk Tone vilde have skrækket Dem. I
1-10 # Sudhher Sandheed jeg Maader ike vover ikke at byke Fulleke dete Mahlae byde Publikum dette Machværk,
1-11 # jeg e er saa forfangelig, forfængelig, at jeg kid, bedd som ei vil træde frem for Vilittm Publikum
1-12 # ndens uden jeg har naged noget at bydl, byde, som forttine Oælsomhad, fortiener Opmærksomhed.
1-13 # sindler Indtil den økd, Tid, som jeg endan e Tu, vaver endnu ei øiner, lever jeg i ukekomdte Ubekjendt
1-14 # Slhed er Stilhed. – Jeg havde haat havt i kende Sinde at skrive til Dig, som kan land
1-15 # mand Deg, mn Tuden i Dag, men Tiden er kalet baret, løbet bort, og jeg maare ogsalte maae opsætte til
1-16 # en anden Gig Gang at gile give Dig begrver mo met landæesom begreber om mit landvæsen,
1-17 # som i Forhald Forhold til den Fødbred haed de Fodbred Jord der er mig givet, kor aveke kan overbe
1-18 # galide vise Dig m, om, at hrdens Tebnde maaskkae Jordens Arbeide maaskee er det, som berde be Bedst
1-19 # Ske mnde Se haænnde hils trives under Diine Hænder. Hils Din Løbe, Lotte.
1-20 # Din alig fengeil evig hengivne
1-21 # be Aall

3.5 Aall - med Norhand - lav epoch:

- 1-1 # ser gave tiltee ser ertei de Grev Moltke i Rentekammeret, som ønsket at kunne
1-2 # har i beges e beoeeeskae, se e je øf faae et begreb om Biergvidenskab ere, som en Green af
1-3 # det seg her n he ogter d Dden ald og ded Fag han nu har opofret sig. Desuden vilde jeg dermed
1-4 # ger mn mste toe e seder e er bedstt give mine norske Venner en Idee om en Videnskab
1-5 # hae hvori de ver illele evidens Kald ted vare aldeeles uvidende. Saaledes torde jeg Dtede aldeeles
1-6 # mige sorlekee hees att, jg mette ved de so ingen Forkundskaber forudsætter, jeg maatte lade dem sma
1-7 # ge nget noget af ill helge bidenkabr, alle Herlige Videnskaberne, og ed e tier ved et kiælent
1-8 # og gentet Jørg lide Dilg de e eed deg de pyntet Sprog lokker Den til at nyde Ting, der i
1-9 # Geseset Die eld hen sædet den d filosofisk Tone vilde have skrækket Dem. I
1-10 # Hdger, vier ed et lde tiled det Mafvr Sandheed jeg vover ikke at byde Publikum dette Machværk,
1-11 # j e sar forforlg, a eg d ld se her berlite jeg er saa forfængelig, at jeg ei vil træde frem for Publikum
1-12 # nede g he mg et ld, se fotem Gyrelbernhed uden jeg har noget at byde, som fortiener Opmærksomhed.
1-13 # Indtil den Tid, som jeg endnu ei øiner, lever jeg i er Dded, ser, ede me øej eleredt Ubekjendt
1-14 # Sfd g haad haa den e fev te deg so ke Stilhed. – Jeg havde havt i Sinde at skrive til Dig, som land
1-15 # med deg mee Diden mand i vives beret, og, me ersadte tel Dag, men Tiden er løbet bort, og jeg maae opsætte til
1-16 # e eden jeg en anden Gang at geve give Dig begrvr ee et londese begreber om mit landvæsen,
1-17 # s dehd som i Forhold til den tidbr hed de eg gevet de md Fodbred Jord der er mig givet, kan overbe
1-18 # gaad ag e, vise Dig om, at hrden vevet moker e d, hen esd Jordens Arbeide maaskee er det, som Bedst
1-19 # se ed va had he D eve trives under Diine Hænder. Hils Din Lotte.
1-20 # D eg feg Din evig hengivne
1-21 # De Aall

Vedlegg 2

Sammenligning av modellenes transkribering av side 71 (Gen-2)

Vedlegg: <https://media.digitalarkivet.no/view/150323/71>

Manuell transkripsjon:

1. Det Korn som Du kan sende mig i Høst v
2. Du søge at faae overtalt *Tellev* eller en ande
3. til at medtake, da det Skib som skal nedtake Ko
4. hos *Madme Cottrup* og *Møller* ei kan rumme
5. meer end jeg der faar. Iligemaade om der gives
6. lighed for det Bryg, som *Møller* vil levere.–
7. Jeg sender Dig med *Tellev* 2^{de} Stukker til Kak
8. lovnen, da jeg ikke erindrer om det var to eller
9. som Mangler, og Du kan da beholde det een
10. i Reserve. Iligemaade sender min Kone Dig en
11. Tylte
12. Tønde Østers og en Dunk ~~Multebær~~.
13. Til Vinteren tager jeg tager jeg fat paa min Afhandli
14. omdanner den til kritisk Bedømmelse; dog vi
- 15.
16. den vist neppe komme for Publikums Øine
17. *Lovise* hilser Din *Lotte* og vi forsikkrer begge at
18. vi erindres vore Venner paa *Hofmannsgave* med ø
19. oprigtigt Venskab.
20. JAall#^z

Sammenligning av de tre etablerte modellene som det bygges på:

Her sammenlignes hver modells transkripsjon av side 71 med den manuelle transkripsjonen. Ord som er farget rødt og krysset ut, er ord modellene har transkribert feil. Ord som er farget grønt er de tilsvarende ordene fra den manuelle transkripsjonen. Ord som er farget hvitt er ord som modellene har transkribert riktig.

NorHand 1820-1940

- 1-1 # Det ~~kav~~ Korn som Du kan ~~hend mig~~ I sende mig i Høst v
1-2 # Du ~~sage~~ søge at faae ~~ødtalt~~ Veller ~~overtalt~~ Tellev eller en ande
1-3 # til at ~~medtage~~, medtakte, da ~~et sil~~ det Skib som skal ~~undtage~~ nedtakte Ko
1-4 # hos ~~modte~~ Madme Cottrup og Møller-Møller ei kan rumme
1-5 # ~~mear~~ meer end jeg der faar. ~~Fligemaade~~ omd Iligemaade om der ~~giker~~ gives
1-6 # ~~ligts~~ lighed for ~~tit bryg~~, det Bryg, som Møller-Møller vil ~~lavee~~ levere.
1-7 # Jeg ~~sende~~ sender Dig med ~~Feller 2te~~ Sypr Tellev 2de Stukker til ~~hod~~ Kak
1-8 # lovnem, da jeg ~~vi~~ ikke erindrer om ~~at~~ det var to ~~aller~~ eller
1-9 # som Du ~~mangler~~, Mangler, og Du kan da beholde det ~~ved~~ een
1-10 # i ~~Leser~~te. Reserve. Iligemaade ~~hande~~ sender min ~~Lone~~ Kone Dig ~~at~~ en
1-11 # ~~jl~~ Tylte
1-12 # ~~Sond~~ Østar Tønde Østers og en ~~Vint~~ muskbær, Dunk Muldebær
1-13 # Til ~~Fintaran~~ taær Vinteren tager jeg ~~tat~~ tager jeg fat paa min ~~Affaen~~ Afhandli
1-14 # ~~omdamner~~ omdanner den til brude, og sender Dig den ~~maad~~ kritisk Bedømmelse; dog vi
1-15 # ~~Blad for~~ Blad til ~~Lritisk~~ bedømmelse; dog d
1-16 # den vist neppe komme for ~~Poblikum~~ Pins-Publikums Øine
1-17 # ~~Lavise~~ Lovise hilser Din ~~lotte~~ Lotte og vi ~~forside~~ forsikkrer begge at
1-18 # vi ~~erindrer~~ som ~~sanner~~ erindres vore Venner paa ~~Kofmannsgave~~ ved a Hofmannsgave med ø
1-19 # ~~oproytist~~ denskab- oprigtigt Venskab.
1-20 # ~~Hallya~~ JAall#

19th century Danish Gothic handwriting

- 1-1 # Det ~~kom~~ Korn som ~~du~~ Du kan sende ~~ingen~~ mig i Høst v
1-2 # ~~de sag~~ Du søge at faae ~~overtalt~~ Tellev eller ~~eller~~ en ~~anden~~ ande
1-3 # til at ~~medtage~~, medtakte, da det Skib ~~ham~~ som skal ~~indtage~~ nedtakte Ko
1-4 # hos ~~Mad~~ Lottrup Madme Cottrup og Møller-Møller ei kan ~~om~~ rumme
1-5 # ~~men~~ ind og de faae meer end jeg der faar. Iligemaade om ~~de~~ ~~gik~~ der gives
1-6 # lighed for det ~~Byg~~, Bryg, som Møller vil ~~lader~~ levere.
1-7 # ~~tog~~ ~~hende~~ dog Jeg sender Dig med ~~eller~~ Tellev 2de ~~stykke~~ Stukker til ~~Kag~~ Kak
1-8 # lovnem, da jeg ~~nu~~ ~~erindre~~ ikke erindrer om det var to eller
1-9 # som ~~de~~ ~~mangle~~, Mangler, og ~~de~~ Du kan da beholde det ~~end~~ een
1-10 # i ~~Risende~~, ligemaade ~~hende~~ Reserve. Iligemaade sender min Kone ~~i~~ Dig en
1-11 # Til Tylte
1-12 # ~~Land~~ Øster Tønde Østers og en ~~de~~ Kirkebar, Dunk Muldebær
1-13 # Til Vinteren ~~tage~~ og ~~tal~~ saa tager jeg tager jeg fat paa min ~~afham~~ Afhandli
1-14 # ~~omdomme~~ omdanner den til ~~bruge~~, og ~~hende~~ det den ~~maa~~ kritisk Bedømmelse; dog vi
1-15 # ~~blad for~~ blod til ~~Litisk~~ Bedømmelse, dog ey
1-16 # den vist neppe komme for ~~Publicum~~ Ons-Publikums Øine
1-17 # ~~have~~ hilser der ~~lotte~~ Lovise hilser Din Lotte og ~~ei~~ 1, 28 aar og vi forsikkrer begge at
1-18 # ~~ei~~ endne var denne en ~~Comunen~~ vi erindres vore Venner paa Hofmannsgave med ~~en~~ ø
1-19 # ~~prigtigt~~ denskab- oprigtigt Venskab.
1-20 # ~~Halle~~ JAall#

18C Danish Administrative Writing (PyLaia)

- 1-1 # ~~dt hver~~ Det Korn som ~~de kand~~ sand mig fik Du kan sende mig i Høst v
- 1-2 # ~~de høge~~ Du søge at faae overtalt Siller-Tellev eller ~~end~~ en ande
- 1-3 # til at medtage, medtakte, da det ~~Rib ham~~ Skib som skal ~~undtage?~~ nedtage Ko
- 1-4 # hos ~~Madm Lottrup~~ Madme Cottrup og ~~Møller~~ Møller ei kan ~~ømme~~ rumme
- 1-5 # ~~men~~ meer end ~~ieg~~ de faae. Iligemaade jeg der faar. Iligemaade om ~~de gike~~ der gives
- 1-6 # lighed for ~~dit byg~~, det Bryg, som ~~Møller ael lever~~ Møller vil levere.
- 1-7 # Jeg ~~hende dig~~ sender Dig med ~~deleur~~ Tellev 2de ~~stykke~~ Stukker til ~~Kad~~ Kak
- 1-8 # lovnen, da jeg ~~ieke erindre~~ ikke erindrer om det var to eller
- 1-9 # som ~~du mangle~~, Mangler, og ~~du han~~ Du kan da beholde det ~~ien~~ een
- 1-10 # i ~~Kehertet~~ Iligemaad ~~hende~~ Reserve. Iligemaade sender min Kone ~~di eg~~ Dig en
- 1-11 # ~~Tylt~~ Tylte
- 1-12 # ~~Tand Østees~~ Tønde Østers og en ~~dun~~ Vuteerbæe Dunk Multebær
- 1-13 # Til Anten tage sig fet har ~~Vinteren~~ tager jeg tager jeg fat paa min ~~Aftende~~ Afhandli
- 1-14 # ~~ømdomme~~ omdanner den til ~~benke~~, og ~~hende dig~~ den maat kritisk Bedømmelse; dog vi
- 1-15 # ~~blad for~~ blad til leiliske bedømmelse; dog et
- 1-16 # den af ~~nepne~~ vist neppe komme for ~~Peblikning~~ Øn Publikums Øine
- 1-17 # ~~Lavin~~ holde din lotte, Lovise hilser Din Lotte og vi ~~foside~~ forsikkrer begge at
- 1-18 # ei vindne van ~~Damme~~ paa ~~Rofmansgare~~ ind mig vi erindres vore Venner paa Hofmannsgave med ø
- 1-19 # ~~oprøgtigt~~ venskab oprigtigt Venskab.
- 1-20 # ~~Hallen~~ JAall#

Sammenligning av 2-Gen modellene:

JAall + Danish Handwriting(130) v0.1

- 1-1 # Det ~~Kon~~Korn som Du kan sende ~~mige~~ Hes-mig i Høst v
1-2 # Du søge at faae overtalt ~~Teller~~Tellev eller en ande
1-3 # til at ~~medtage~~, medtake, da det Skib som skal ~~indtage~~ ~~ti~~ nedtake Ko
1-4 # hos Madme Cottrup og Møller ei kan ~~rumme~~ rumme
1-5 # meer end jeg der faar. Iligemaade om der gives
1-6 # lighed for det ~~Byg~~, Bryg, som Møller vil ~~lavere~~ levere.
1-7 # Jeg sender Dig med ~~deller~~ Tellev 2de ~~Stykker~~ Stukker til ~~Ka~~ Kak
1-8 # lovnen, da jeg ikke erindrer om det var ~~te~~ to eller
1-9 # som ~~Du mangler~~, Mangler, og Du kan da beholde det een
1-10 # i ~~Bidherves~~ Sligemaade Reserve. Iligemaade sender min Kone Dig ~~ei~~ en
1-11 # ~~Tylke~~ Tylte
1-12 # Tønde Østers og en ~~Dunke~~ Mukterbær.—Dunk Multebær
1-13 # Til Vinteren tager jeg ~~tager~~ jeg fat paa min ~~Affangt~~ Afhandli
1-14 # ~~omdomne~~ omdanner den til Breve, og sender Dig ~~den~~ ~~maat~~ kritisk Bedømmelse; dog vi
1-15 # ~~blad for blod~~ til ~~Kritisk~~ bedømmelse; dog ig ...
1-16 # den vist neppe komme for ~~Pablikum~~ Ping-Publikums Øine
1-17 # Lovise hilser Din Lotte og vi ~~forsikkre~~ begtge og ~~forsikkre~~ begge at
1-18 # vi ~~erindrer~~ ~~vare~~ erindres vore Venner paa Hofmannsgave med ~~ang~~ ø
1-19 # oprigtigt Venskab.
1-20 # ~~JAall~~ JAall#

JAall + Danish Handwriting(150) v0.1

- 1-1 # ~~Dit Kone~~ Det Korn som Du ~~laar~~ kan sende ~~mige~~ Has-mig i Høst v
1-2 # Du søge at faae overtalt ~~Tiller~~ Tellev eller en ande
1-3 # til at ~~medtage~~, medtake, da det Skib som skal ~~undtage~~ ~~Ke~~ nedtake Ko
1-4 # hos Madme Cottrup og Møller ei kan rumme
1-5 # ~~mee~~ meer end jeg der faar. Iligemaade om der gives
1-6 # lighed for det ~~Byg~~, Bryg, som Møller vil levere.
1-7 # Jeg sender Dig med ~~Teller~~ Tellev 2de ~~Stykker~~ Stukker til ~~Kor~~ Kak
1-8 # lovnen, da jeg ikke erindrer om det var ~~ti~~ to eller
1-9 # som ~~Du mangle~~, Mangler, og Du kan da beholde det een
1-10 # i ~~Riserves~~ Reserve. Iligemaade ~~hender~~ sender min Kone Dig ~~ei~~ en
1-11 # ~~Tylt~~ Tylte
1-12 # Tønde ~~Østere~~ Østers og en ~~Dinke~~ Mukterbær.—Dunk Multebær
1-13 # ~~Tile~~ Til Vinteren tager jeg ~~tager~~ jeg fat paa min ~~Affant~~ Afhandli
1-14 # ~~omdomne~~ omdanner den til Breve, og sender Dig ~~den~~ ~~maat~~ kritisk Bedømmelse; dog vi
1-15 # ~~blad for blod~~ til ~~kritisk~~ bedømmelse; dog iilt ...
1-16 # den vist neppe komme for ~~Pablikune~~ ~~Øin~~ Publikums Øine
1-17 # Lovise hilser Din Lotte og vi ~~forsikkre~~ forsikkre begge ~~og~~ at
1-18 # vi ~~erindre~~ ~~vare~~ erindres vore Venner paa Hofmannsgave med ~~ang~~ ø
1-19 # ~~oprotigt~~ venskab. oprigtigt Venskab.
1-20 # ~~JAall~~ JAall#

JAall + Norhand stor v.02

- 1-1 # ~~Dail~~ Kare-~~Det~~ Korn som ~~Dku~~-Du kan ~~sene~~ mngge. ~~S9~~-sende mig i Høst v
1-2 # ~~Dun~~-Du søge at faae overtalt ~~Killev~~-Tellev eller en ~~øndt~~-ande
1-3 # til at ~~medtage~~-medtake, da det Skib som skal ~~indtage~~ til-~~nedtake~~ Ko
1-4 # ~~hers~~ Modmei-~~Cottrug~~-hos Madme Cottrup og Møller ei kan ~~rummin~~-rumme
1-5 # ~~men~~-meer end jeg der faar. ~~Sligemaade~~-Iligemaade om der ~~giver~~-gives
1-6 # lighed for det ~~Byg~~-Bryg, som Møller-Møller vil ~~levere~~-levere.
1-7 # Jeg sender Dig med ~~Tiller~~-vide ~~Styker~~-Tellev 2de Stukker til ~~Kiar~~-Kak
1-8 # lovnen, da jeg ikke erindrer om det var ~~to~~-to eller
1-9 # som ~~Du~~ mangler, Mangler, og Du kan da beholde det ~~een~~-een
1-10 # i ~~iherves~~ Kigemaarde Reserve. Iligemaade sender min Kone Dig ~~ei~~-en
1-11 # Tylte
1-12 # Tønde ~~Østere~~-Østers og en ~~Dunke~~ Mntrbær-~~Dunk~~ Multebær
1-13 # Til Vinteren ~~tiger~~-tager jeg ~~hlis~~-tager jeg fat paa min ~~Afmas~~-Afhandli
1-14 # ~~ømdommer~~-ømdanner den til ~~breve~~, og sender Dig den molt kritisk Bedømmelse; dog vi
1-15 # ~~blad~~ for ~~blad~~ til kritisk Bedømmelse, deg vet...
1-16 # den vist ~~nepege~~-neppe komme for ~~Publiken~~ Øin-~~Publikums~~ Øine
1-17 # Lovise hilser Din Lotte og vi forsikkrer begge ~~det~~-at
1-18 # vi ~~erendrer~~ voare-erindres vore Venner paa ~~Hofmannssans~~-Hofmannsgave med ~~øis~~-ø
1-19 # ~~øprygtigs~~ Venstag-~~øprigtigt~~ Venskab.
1-20 # JAall#-JAall#

JAall - 150 - v0.1

- 1-1 # ~~D~~kor-~~Det~~ Korn som ~~Dli~~-Du kan ~~see~~-sende mig ~~l~~-i Høst v
1-2 # Du ~~sig~~-søge at ~~saae~~ ovetalt ~~Sikke~~ elle-faae overtalt Tellev eller en ~~øndt~~-ande
1-3 # til at ~~medtage~~-medtake, da det ~~Skil~~-sei Skib som skal ~~indtage~~-nedtake Ko
1-4 # hos ~~Mode~~ Cottrag-Madme Cottrup og ~~Mølle~~-Møller ei kan ~~ummin~~-rumme
1-5 # ~~mee~~-meer end jeg ~~de~~ saae. ~~Slegemaade~~ der faar. Iligemaade om ~~de~~ give-der gives
1-6 # ~~lighe~~ lighed for det ~~Byg~~-soa Møle-Bryg, som Møller vil ~~levee~~-levere.
1-7 # Jeg ~~sende~~-sender Dig med ~~Tele~~-Tellev 2de ~~Stykke~~-Stukker til ~~al~~-Kak
1-8 # lovnen, da jeg ikke erindrer om det var ~~ta~~-elle-to eller
1-9 # ~~sam~~-Du mangle, som Mangler, og Du kan da beholde-beholde det ~~iii~~-een
1-10 # i ~~Sherde~~ Rligemaar ~~sende~~ man ~~Koae~~ Di-eg-Reserve. Iligemaade sender min Kone Dig en
1-11 # Tylt Tylte
1-12 # Tønd ~~Østere~~-Tønde Østers og en ~~Dunt~~ Ontebe-~~Dunk~~ Multebær
1-13 # ~~Al~~ Vinteen tage-Til Vinteren tager jeg ~~sot~~-tager jeg fat paa ~~min~~ Gsaa-min Afhandli
1-14 # ~~ømdome~~-ømdanner den til ~~breve~~, og ~~sende~~ Dig den ~~mæl~~-kritisk Bedømmelse; dog vi
1-15 # ~~blad~~ for ~~blad~~ til kritisk ~~Bdømmmelse~~-d0 Ot ...
1-16 # ~~dea~~ Alst ~~nepe~~-den vist neppe komme for ~~Pøbliken~~ i-~~Publikums~~ Øine
1-17 # Lovise ~~hilse~~-hilser Din Lotte og vi ~~fesikte~~-forsikkrer begge at
1-18 # vi ~~erndre~~ van ~~Sanne~~-erindres vore Venner paa ~~Hofmasnspaae~~ mad ~~li~~-Hofmannsgave med ø
1-19 # ~~øpta~~ Vanskab-~~øprigtigt~~ Venskab.
1-20 # Jaall#-JAall#

JAall + 18C Danish Administrative Writing v0.1

- 1-1 # Det ~~vore~~ Korn som ~~ike~~ Du kan ~~sens~~ sende mig ~~his~~ i Høst v
1-2 # Du søge at faae overtalt ~~Tiller~~ Tellev eller en ande
1-3 # til at ~~medtage~~, ~~medtake~~, da det Skib som skal ~~indtage~~ til ~~nedtake~~ Ko
1-4 # hos ~~Modme~~ Lottrup ~~Madme~~ Cottrup og Møller ei kan ~~ømmin~~ rumme
1-5 # meer ~~md~~ end jeg ~~de~~ der faar. Iligemaade om ~~de~~ der gives
1-6 # lighed for det ~~byg~~, ~~Bryg~~, som ~~Møller~~ Møller vil ~~leve~~ levere.
1-7 # Jeg sender Dig med ~~dellev~~ Tellev 2de ~~Stykke~~ Stukker til ~~Kot~~ Kak
1-8 # lovnen, da jeg ikke erindrer om det var to eller
1-9 # som ~~Du~~ mangler, ~~Mangler~~, og Du kan da beholde det een
1-10 # i ~~Kiherdes~~ Reserve. Iligemaade sender min Kone ~~Di~~ ei Dig en
1-11 # Tylte
1-12 # ~~Tand~~ Tønde Østers og en ~~Dine~~ Mutterbær, ~~Dunk~~ Muldebær
1-13 # Til Vinteren tager jeg ~~far~~ tager jeg fat paa min ~~Afsand~~ Afhandli
1-14 # ~~omdomme~~ omdanner den til breve, og sender Dig ~~den~~ maate kritisk Bedømmelse; dog vi
1-15 # ~~blad~~ for ~~blod~~ til ~~kriviske~~ bedømmelse, dog vil ...
1-16 # den vist ~~neppe~~ neppe komme for ~~Poblikune~~ sine ~~Publikums~~ Øine
1-17 # Lovise ~~helser~~ hilser Din ~~lotte~~ Lotte og vi ~~forsikkre~~ forsikkrer begge at
1-18 # vi ~~erindrer~~ vare ~~Vanner~~ erindres vore Venner paa ~~Hofmanns~~ gave ~~Hofmannsgave~~ med ~~oug~~ ø
1-19 # ~~oprygtigt~~ venskab, oprigtigt Venskab.
1-20 # JAall# JAall#

JAall + Norhand split v0.1

- 1-1 # Uls kod skor dide lert stygd ptegl. Det Korn som Du kan sende mig i Høst v
1-2 # ~~Hislg~~ Du søge at ~~save~~ elrtedti ~~Pifkle~~ iter jeg, ~~ornde~~ faae overtalt Tellev eller en ande
1-3 # ~~sil~~ til at ~~mydtage~~, ~~medtake~~, da det ~~Pkil~~ sai Skib som skal ~~yndtags~~ to ~~nedtake~~ Ko
1-4 # hos ~~Mdde~~ Cottrup ~~Madme~~ Cottrup og Møller ~~vi~~ ei kan ~~øemi~~ rumme
1-5 # ~~mae~~ jeedt sjeeg meer end jeg ~~der~~ ~~fabeds~~ Sligemaade, faar. Iligemaade om der ~~jjøeds~~ gives
1-6 # lighed for det ~~Byg~~, ~~Bryg~~, som ~~Malles~~ Møller vil ~~laved~~ levere.
1-7 # ~~de~~ sande Jeg sender Dig med ~~dkeller~~ 5de ~~Styde~~ Pølstøts ~~Tellev~~ 2de Stukker til Kak
1-8 # ~~lavnen~~, lovnen, da jeg ikke erindrer ~~ome~~ om det var to ~~elte~~ eller
1-9 # ~~saa~~ Du maagte, som ~~Mangler~~, og Du kan da ~~Pfoldte~~ beholde det ~~til~~, een
1-10 # ~~ikkeves~~ Slyenaad skiodes ~~mn~~ Pove. 1 eg i Reserve. Iligemaade sender min Kone Dig en
1-11 # ~~Tle~~ Tylte
1-12 # ~~Sande~~ sterr Tønde Østers og en ~~Dinte~~ Muktekde, ~~Dunk~~ Muldebær
1-13 # ~~ild~~ Paberes tige Et ~~sit~~ ska my ~~Atls~~ Til Vinteren tager jeg tager jeg fat paa min Afhandli
1-14 # ~~oidommer~~ dan omdanner den til breve, og sender Dig ~~dan~~ ngt. kritisk Bedømmelse; dog vi
1-15 # ~~fkld~~ for ~~Slad~~ til ~~kodtisk~~ Bdømmelse 1 tt ...
1-16 # ~~davs~~ dist vepe kame den vist neppe komme for ~~Poblikunns~~ 10d. ~~Publikums~~ Øine
1-17 # ~~lavis~~ hilse Di ~~latt~~, Lovise hilser Din Lotte og vi ~~feskka~~ beggedet. forsikkrer begge at
1-18 # vi ~~eruider~~ vaae ~~skrive~~ erindres vore Venner paa ~~Pofasenshavs~~ md ~~hgst~~ Hofmannsgave med ø
1-19 # ~~optg~~ Vanskal, oprigtigt Venskab.
1-20 # Jsallse JAall#