

Resum

Actualment existeixen grans col·leccions de documents manuscrits en llibreries de tot el món. La gran demanda d'aquests recursos ha portat a la creació de llibreries digitals per tal de facilitar la preservació i accés electrònic a aquests documents. No obstant, la transcripció de les imatges d'aquests documents no està sempre disponible per tal de permetre una cerca ràpida i eficaç als usuaris, o d'extraure patrons i dades estadístiques automàticament. Aquesta tesi presenta una nova aproximació per a la transcripció assistida per ordinador (CAT) de documents de text manuscrit emprant sistemes de reconeixement de text manuscrit (HTR).

L'objectiu de les aproximacions CAT és, completar de manera eficaç una tasca de transcripció mitjançant la col·laboració home-màquina, ja que l'esforç requerit per generar una transcripció manual és alt, i les transcripcions obtingudes automàticament per sistemes estat del art encara no arriben a la precisió requerida. Aquesta tesi es centra en una aplicació especial de CAT, que és la transcripció de documents manuscrits antics quan l'esforç d'usuari és limitat, i en conseqüència, el document no pot ser revisat completament. En aquesta aproximació, l'objectiu és generar la millor transcripció possible emprant l'esforç d'usuari disponible. Aquesta tesi ofereix una guia completa del procés de CAT des de l'extracció de característiques fins a l'interacció d'usuari.

Primer, es proposa una aproximació estadística per a generalitzar la transcripció interactiva. Donat que la seua aplicació directa és inabordable, s'han realitzat una sèrie d'assumpcions per tal d'aplicar-la en dos tasques diferents: la transcripció interactiva de documents de textos manuscrits i la detecció del format de documents de text.

A continuació, es descriu el procés de digitalització i anotació de dos documents manuscrits antics reals. Aquest procés s'ha portat a terme donat el nombre escàs de recursos similars i la necessitat de dades anotades per tal de comprovar totes les eines i tècniques desenvolupades en aquesta tesi. Aquests dos documents han estat escollits amb cura amb l'objectiu de representar les típiques dificultats que es troben al utilitzar tècniques HTR. Es presenten resultats de referència en aquests dos documents obtinguts amb un sistema estàndard per tal de servir de referència. Finalment, aquests documents s'han fet públics i accessibles lliurement a la comunitat. Hi ha de tindre en compte que totes les tècniques i mètodes desenvolupats en aquesta tesi s'han evaluat en aquests dos documents antics.

Seguidament, s'estudia i verifica de manera exhaustiva una aproximació CAT per HTR quan l'esforç d'usuari és limitat. L'objectiu final d'aplicar CAT s'aconsegueix mitjançant l'unió de tres processos separats. Donat el reconeixement automàtic d'un sistema HTR. El primer procés consisteix en localitzar paraules (possiblement) incorrectes i emprar l'esforç d'usuari disponible en supervisar-les i corregir-les (si és necessari). Donat que la majoria de les paraules no es van a supervisar ja que sols hi ha una quantitat limitada d'esforç d'usuari, sols unes poques sern seleccionades per una estimació de la seua correctitud, o per a ser més precises, seleccionades d'acord amb el seu nivell de confiança. A continuació, el segon procés comença una vegada aquestes paraules de baixa confiança han estat revisades. Aquest procés actualitza el reconeixement del document tenint en compte les cor-

recions, el qual millora la qualitat de les paraules que no han estat revisades per l'usuari. Finalment, l'ltim procs adapta el sistema a partir de l'ltima transcripci parcialment supervisada (i possiblement imperfecta) que s'ha obts. En aquesta adaptaci, el sistema escolleix de manera intelligent que paraules correctes de la transcripci son utilitzades en l'adaptaci. En consequncia, el sistema adaptat reconeixer millor les futures transcripcions. Els experiments de transcripci realitzats utilitzant aquesta aproximaci CAT mostren que aquesta aproximaci es ms efica quan l'esfor d'usuari aplicat es baix.

L'ltima contribuci d'aquesta tesi es un mtode per a equilibrar la qualitat de transcripci final i l'esfor de supervisi aplicat quan s'utilitza l'aproximaci CAT previament descrita. En altres paraules, aquest mtode permeteix al usuari controlar la quantitat d'errors en les transcripcions obtses amb una aproximaci CAT. La motivaci d'aquest mtode es permetre als usuaris decidir la qualitat final desitjada en els document, ja que una transcripci parcialment errnia pot ser sufficient per a entendre el contingut, i l'esfor requerit per obtindre aquesta transcripci pot ser significativament menor que el d'obtindre una transcripci manual completa. Com a resultat, el sistema estima l'esfor d'usuari mnim requerit per alcanar la quantitat d'error definit pel usuari. L'estimaci del error es realitza calculant per separat l'error causat per cada paraula reconeguda, per a desprs demanar al usuari que revis aquelles on hi ha ms errors.

A ms, es presenta un prototip interactiu que integra la majoria de les tcniques interactives presentades en aquesta tesi. Aquest prototip s'ha desenvolupat per a ser utilitzat per experts paleogrfsics, que no poseixen cap coneiximent de les tecnologies HTR. Desprs de ser ajustats per experts en HTR, el prototip permet als transcriptors anotar un document manualment o utilitzar l'aproximaci CAT presentada. Tots els processos automtics, com el reconeixement, s'executen en segn pla abstraent al transcriptor dels detalls interns del sistema. El prototip va ser probat per un expert transcriptor i desmostr ser adequat i eficient per a la seua finalitat. El prototip est disponible lliure i publicament mitjanant una llicencia GNU (GPL).