

Parallele Korpora

Universität München, CIS

“Korpus- und UNIX-Tools”
Sommersemester 2006

Dozenten: Sebastian Nagel, Yeong Su Lee

Natalya Shupletsova, Stefan Partusch

Was ist ein paralleler Korpus?

Sammlung von **identischen** Texten
in **verschiedenen** Sprachen:

- zweisprachig (bilingual)
- mehrsprachig (multilingual)

Was ist ein paralleler Korpus?

Die Idee einen Text parallel in mehreren Sprachen anzubieten ist nicht neu.

A Scifi. 46.a. Hec noia filioꝝ israel: qui in greſſi ſunt in egyptu cū iacob pfe eoz: vnus porouneon eis aiyuſonima iacob w τω περιουτων. ε qisq cū domibꝫ ſuis introierūt: rubē: ſymeō: leuī: καρος πανοικι αυτων εισηλθουσα, ρουβην, συμειων, λευι iudas: ifachar: zabulo: τ beniamin: da: et neph ioudas, ισαχαρ, αβουλων, και βενιαμιν, δαν, καινεφ θηλι: gad: τ aser. ioseph at erat in egypto. erat at θαλι, γαδ, και ασηρ. ιωσηφ δε ην εν αιγυπτω. ησασ δε οες sic q egressē fūt ex iacob: qnq τ sep τασαι ψυχαι αι εξελευσαι εε ιακωβ, τεντε και εβ tuaginta. mortuus eāt ioseph τ οες fra= do μηκοντα. ετελευτησε δε ιωσηφ και παντες οι αδελ tres eius: τ οis għatto illa. at filij israel cre= φοι αυτ, και πασα η γενεα εκεινη. οι δε υιοι ισραηλ ηυξη uerūt: τ multiplicati fūt: τ abidātes fuerūt: τ in= θισασ, και επι πληθυνθισασ, και χυδ αι οι εθνοντο, και κα ualuerūt valde multiplicati. multiplicauit at fra illos. sur τ ιχουον σφοδρα σφοδρα. επληθυνε δε η γη αυτ. ανε rexit at rex alter sup egyptu: q nō cognoscebat ση δε βασιλευς ετερος εω αυγυπτον, ος ουκ ηδει τον ioseph. dixit at genti sue: ecce gens fi iωσιφ. ειπαε δε τω εθνε αυτου, ιδου γο εθνος των υι lioz israel magna valde multitudo: τ pualet sup nos. ων ισραηλ μεγα πολυ πληθος, και ιχυει υπερ η= venite ergo sapiēter opprimatꝫ eos: ne forte μάς. δευτε ουν κατασοφισμεθα αυτους, μη ποτε multiplieēt: τ qñ acciderit nobis bellū: ad= πληθυνθη, και ηνικα αν συμβη ημιν πολεμος, προ= detur τ isti ad aduersarios. τ δε στεθησον εμ και αυτοι προς τους υπεναντιους. και εκ bellantes nos: egrediētur de fra. τ pre= πολεμισαστες ημας, εξελευσον εμ εκ της γης. και εωε fecit eis pfectos operū: vt affligerēt σησεν αυτοις επιστατας των εργαων, ινα κακωσωσιν eos in operibꝫ. τ edificaueēt ciuitates munitas αυτ. εν τοις εργαοις. και ακοδ ομοσασ πολεις θυρατας pharaōi: τ phihth et ramesse: τ on: q eit τω φαραω, την τε φιδωμ και ραμεση, και ων, η εστι η hellopolis. quāto at eos humiliabāt: tāto λιουπολις. και οτι δε αυτους εταπεινον, τοσούτω plures fiebant. τ inualuerūt valde. τ abhominatio πλειους εγινοντο. και ιχυον σφοδρα. και εβδελυσον ne huerūt egyptij a filijs israel. τ oppresse= το οι αιγυπτιοι ακωε των υιων ισραηλ. και κατεδυνα= rūt egyptij filios israel vi: τ afflix= σευον οι αιγυπτιοι τς υιους ισραηλ βια, και κατωδ υ= rūt eoz vitā in operibꝫ duris in των αυτων την ζωην εν τοις εργαοις τοις σκληροις, εν τω into τ lateritio: τ oibꝫ operibꝫ q in a= πικλω και τη laterithia, και πασι τοις εργαοις τοις εν τοις πε gris: fm oia opera: qbꝫ in fuit ut redegeēt eos διοις, και τα πάντα τα εργα, ων κατεδουλουντο αυτους

B Inter p. cl. al. Hec noia filioꝝ israel: qui in greſſi ſunt in egyptu cū iacob pfe eoz: vnus porouneon eis aiyuſonima iacob w τω περιουτων. ε qisq cū domibꝫ ſuis introierūt: rubē: ſymeō: leuī: καρος πανοικι αυτων εισηλθουσα, ρουβην, συμειων, λευι iudas: ifachar: zabulo: τ beniamin: da: et neph ioudas, ισαχαρ, αβουλων, και βενιαμιν, δαν, καινεφ θηλι: gad: τ aser. ioseph at erat in egypto. erat at θαλι, γαδ, και ασηρ. ιωσηφ δε ην εν αιγυπτω. ησασ δε οες sic q egressē fūt ex iacob: qnq τ sep τασαι ψυχαι αι εξελευσαι εε ιακωβ, τεντε και εβ tuaginta. mortuus eāt ioseph τ οες fra= do μηκοντα. ετελευτησε δε ιωσηφ και παντες οι αδελ tres eius: τ οis għatto illa. at filij israel cre= φοι αυτ, και πασα η γενεα εκεινη. οι δε υιοι ισραηλ ηυξη uerūt: τ multiplicati fūt: τ abidātes fuerūt: τ in= θισασ, και επι πληθυνθισασ, και χυδ αι οι εθνοντο, και κα ualuerūt valde multiplicati. multiplicauit at fra illos. sur τ ιχουον σφοδρα σφοδρα. επληθυνε δε η γη αυτ. ανε rexit at rex alter sup egyptu: q nō cognoscebat ση δε βασιλευς ετερος εω αυγυπτον, ος ουκ ηδει τον ioseph. dixit at genti sue: ecce gens fi iωσιφ. ειπαε δε τω εθνε αυτου, ιδου γο εθνος των υι lioz israel magna valde multitudo: τ pualet sup nos. ων ισραηλ μεγα πολυ πληθος, και ιχυει υπερ η= venite ergo sapiēter opprimatꝫ eos: ne forte μάς. δευτε ουν κατασοφισμεθα αυτους, μη ποτε multiplieēt: τ qñ acciderit nobis bellū: ad= πληθυνθη, και ηνικα αν συμβη ημιν πολεμος, προ= detur τ isti ad aduersarios. τ δε στεθησον εμ και αυτοι προς τους υπεναντιους. και εκ bellantes nos: egrediētur de fra. τ pre= πολεμισαστες ημας, εξελευσον εμ εκ της γης. και εωε fecit eis pfectos operū: vt affligerēt σησεν αυτοις επιστατας των εργαων, ινα κακωσωσιν eos in operibꝫ. τ edificaueēt ciuitates munitas αυτ. εν τοις εργαοις. και ακοδ ομοσασ πολεις θυρατας pharaōi: τ phihth et ramesse: τ on: q eit τω φαραω, την τε φιδωμ και ραμεση, και ων, η εστι η hellopolis. quāto at eos humiliabāt: tāto λιουπολις. και οτι δε αυτους εταπεινον, τοσούτω plures fiebant. τ inualuerūt valde. τ abhominatio πλειους εγινοντο. και ιχυον σφοδρα. και εβδελυσον ne huerūt egyptij a filijs israel. τ oppresse= το οι αιγυπτιοι ακωε των υιων ισραηλ. και κατεδυνα= rūt egyptij filios israel vi: τ afflix= σευον οι αιγυπτιοι τς υιους ισραηλ βια, και κατωδ υ= rūt eoz vitā in operibꝫ duris in των αυτων την ζωην εν τοις εργαοις τοις σκληροις, εν τω into τ lateritio: τ oibꝫ operibꝫ q in a= πικλω και τη laterithia, και πασι τοις εργαοις τοις εν τοις πε gris: fm oia opera: qbꝫ in fuit ut redegeēt eos διοις, και τα πάντα τα εργα, ων κατεδουλουντο αυτους

Inter p. cl. al. Hec noia filioꝝ israel: qui in greſſi ſunt in egyptu cū iacob pfe eoz: vnus porouneon eis aiyuſonima iacob w τω περιουτων. ε qisq cū domibꝫ ſuis introierūt: rubē: ſymeō: leuī: καρος πανοικι αυτων εισηλθουσα, ρουβην, συμειων, λευι iudas: ifachar: zabulo: τ beniamin: da: et neph ioudas, ισαχαρ, αβουλων, και βενιαμιν, δαν, καινεφ θηλι: gad: τ aser. ioseph at erat in egypto. erat at θαλι, γαδ, και ασηρ. ιωσηφ δε ην εν αιγυπτω. ησασ δε οες sic q egressē fūt ex iacob: qnq τ sep τασαι ψυχαι αι εξελευσαι εε ιακωβ, τεντε και εβ tuaginta. mortuus eāt ioseph τ οες fra= do μηκοντα. ετελευτησε δε ιωσηφ και παντες οι αδελ tres eius: τ οis għatto illa. at filij israel cre= φοι αυτ, και πασα η γενεα εκεινη. οι δε υιοι ισραηλ ηυξη uerūt: τ multiplicati fūt: τ abidātes fuerūt: τ in= θισασ, και επι πληθυνθισασ, και χυδ αι οι εθνοντο, και κα ualuerūt valde multiplicati. multiplicauit at fra illos. sur τ ιχουον σφοδρα σφοδρα. επληθυνε δε η γη αυτ. ανε rexit at rex alter sup egyptu: q nō cognoscebat ση δε βασιλευς ετερος εω αυγυπτον, ος ουκ ηδει τον ioseph. dixit at genti sue: ecce gens fi iωσιφ. ειπαε δε τω εθνε αυτου, ιδου γο εθνος των υι lioz israel magna valde multitudo: τ pualet sup nos. ων ισραηλ μεγα πολυ πληθος, και ιχυει υπερ η= venite ergo sapiēter opprimatꝫ eos: ne forte μάς. δευτε ουν κατασοφισμεθα αυτους, μη ποτε multiplieēt: τ qñ acciderit nobis bellū: ad= πληθυνθη, και ηνικα αν συμβη ημιν πολεμος, προ= detur τ isti ad aduersarios. τ δε στεθησον εμ και αυτοι προς τους υπεναντιους. και εκ bellantes nos: egrediētur de fra. τ pre= πολεμισαστες ημας, εξελευσον εμ εκ της γης. και εωε fecit eis pfectos operū: vt affligerēt σησεν αυτοις επιστατας των εργαων, ινα κακωσωσιν eos in operibꝫ. τ edificaueēt ciuitates munitas αυτ. εν τοις εργαοις. και ακοδ ομοσασ πολεις θυρατας pharaōi: τ phihth et ramesse: τ on: q eit τω φαραω, την τε φιδωμ και ραμεση, και ων, η εστι η hellopolis. quāto at eos humiliabāt: tāto λιουπολις. και οτι δε αυτους εταπεινον, τοσούτω plures fiebant. τ inualuerūt valde. τ abhominatio πλειους εγινοντο. και ιχυον σφοδρα. και εβδελυσον ne huerūt egyptij a filijs israel. τ oppresse= το οι αιγυπτιοι ακωε των υιων ισραηλ. και κατεδυνα= rūt egyptij filios israel vi: τ afflix= σευον οι αιγυπτιοι τς υιους ισραηλ βια, και κατωδ υ= rūt eoz vitā in operibꝫ duris in των αυτων την ζωην εν τοις εργαοις τοις σκληροις, εν τω into τ lateritio: τ oibꝫ operibꝫ q in a= πικλω και τη laterithia, και πασι τοις εργαοις τοις εν τοις πε gris: fm oia opera: qbꝫ in fuit ut redegeēt eos διοις, και τα πάντα τα εργα, ων κατεδουλουντο αυτους

Incipit liber hellesmoth que nos Exodi dicimus. Hec sūt noia fi .ca. i. filiorū israel: qui in greſſi ſunt i egyptū cū iacob. Singuli cū domibus ſuis introierunt. Rubē: Symeon: Leuī Iudas: ifachar: zabulo & Beniamin. Dā Nephthalim. Gad: & Aser. Erāt igitur oēs aie corū q egressi ſunt de femore iacob septuaginta. Ioseph at in egypto erat: Quo mortuo & vniversis fribus eius oisq cognatiōe ſua: filii israel creuerūt: & quasi germinates multiplicati ſunt: ac roborati nimis impleuerūt terram. Surrexit iterea rex novus super egyptum: qui ignorabat ioseph. & ait ad populū ſuū. Ecce pplis filiorū isrl multus & fortior nobis est. Venite sapiēter opprimamus eū: ne forte multiplicetur: & ſi in gruerit cōtra nos bellū: ad datur inimicis nrīs: expugnatiſq nobis egre diatur de terra. Prepo ſuit itaq: eis magiſtros opuz: vt affligerēt eos oneribus. Edificauerūt q: vrbes tabernaculoꝝ rum pharaōi phiton & ramesse. quātoq op primebāt eos: tāto magis multiplicabātur: & creſcebat. O derātq: filios israel egyptij: & affligēbant illudētes & inuidentes eis: atq ad amaritudinem perdūcebant vitam corū: operibus duris luti ooo & lateris: p omniq: fa= mulatu oooooooooooooo & lateris: p omniq: fa= mulatu oooooooooooooo quo in terre operibus ooooooooooooooooooooo

Transla. Chal. Incipit liber Exodus. Cap. i. Hec sunt nomina filioꝝ israel: qui ingressi sunt in egyptū cum iacob: singuli cum viris domus sue introierūt. Ruben: symeon: leui: et iudas: isachar: zabulon: et beniamin: dan: et nephtalim: et aser. Erant autem aie egresserunt de femore iacob septuaginta. Ioseph autem in egypto erat: et post mortem eius et omnibus fratribus eius: et cognatione sua: filii israel creuerunt: et quasi germinantes multiplicati sunt: et roborati sunt nimis: et impleverunt terram. Surrexit autem post haec rex novus super egyptum: qui nō cognoscebat ioseph. et ait ad populum suum: Ecce populus filiorum israel multus est et fortior nobis est. Venite ergo sapienter opprimamus eum: ne forte multiplicetur: et si in gruerit contra nos bellum: ad datur inimicis nostris: expugnatiſque nobis egre diatur de terra. Prepositi sunt itaque eis magistros operum: ut affligerent eos oneribus. Edificaverunt quoque: vrbes tabernaculoꝝ rum pharaōi phiton & ramesse. quātoque op primebāt eos: tāto magis multiplicabātur: & creſcebat. O derātq: filios israel egyptij: & affligēbant illudētes & inuidentes eis: atq ad amaritudinem perdūcebant vitam corū: operibus duris luti ooo & lateris: p omniq: fa= mulatu oooooooooooooo & lateris: p omniq: fa= mulatu oooooooooooooo quo in terre operibus ooooooooooooooooooooo

וְאֵלֶּה שְׁמוֹת בְּנֵי יִשְׂרָאֵל הַבָּאִים מִצְרַיִם אִישׁ וּבֵיתוֹ בָּאוּ: רְאוּבֵן שִׁמְעוֹן לֵוִי וַיהוּדָה: יִשְׁשַׁכָּר זְבוּלֹן וּבְנִימֵן: דָּן וְנַפְתָּלִי גַד וְאָשֶׁר: וַיְהִי כֹל נַפְשׁ יִצְאֵי יִרְךָ יַעֲקֹב שְׁבַע עִים נַפְשׁ וַיּוֹסֶף הִיָּה בְּמִצְרַיִם: וַיָּמָת יוֹסֵף וְכָל אָחָיו וְכָל הַדּוֹר הַהוּא: וּבְנֵי יִשְׂרָאֵל פָּרוּ וַיִּשְׂרְצוּ וַיִּרְבּוּ וַיַּעֲצֻמוּ כַּמְאֹד כַּמְאֹד: וַתִּמְלֵא הָאָרֶץ אֹתָם: וַיִּקַּם מֶלֶךְ חָדָשׁ עַל מִצְרַיִם אִשֵּׁר לֹא יָדַע אֶת יוֹסֵף: וַיֹּאמֶר אֶל עַמּוֹ הֲנֵה עַם בְּנֵי יִשְׂרָאֵל רַב וְעָצוּם מִמֶּנּוּ: הֲבֵנָה נִתְחַכְמָה לוֹ: כִּן יִרְבֶּה וַיְהִי כִּי תִקְרָאנָה מִלְחָמָה וְנוֹסֶף נָם הוּא עַל שְׂנְאֵינוּ וְנִלְחַם בָּנוּ: וְעַל מֶן הָאָרֶץ: וַיִּשְׁיִמוּ עָלָיו שָׂרִי מִסִּים לַמַּעַן עַנְתּוֹ: בְּסַבְלָתָם וַיִּבְּן עָרֵי מִסְכְּנוֹת לַפְּרֵעָה: אֶת פְּתָם וְאֶת רַעַמְסֵם: וַיִּקְרָא אֶת שְׂרָפִיטֹן וְרַעַמְסֵם: וַיִּבְּנוּ מִכְּנֵי בְנֵי יִשְׂרָאֵל: וַיַּעֲבְדוּ מִצְרַיִם אֶת בְּנֵי יִשְׂרָאֵל: בְּפִרְךָ: וַיִּמְרְרוּ אֶת חַיֵּיהֶם בְּעַבְדָּה קָשָׁה: בְּחָמֶר וּבְלִבְנִים: וּבְכָל עַבְדָּה בְשָׂדֶה: אֶת כָּל עַנְדָּתָם: אֶשֶׁר עָבְדוּ בָהֶם

Scifi. 46.a. Hec noia filioꝝ israel: qui in greſſi ſunt in egyptu cū iacob pfe eoz: vnus porouneon eis aiyuſonima iacob w τω περιουτων. ε qisq cū domibꝫ ſuis introierūt: rubē: ſymeō: leuī: καρος πανοικι αυτων εισηλθουσα, ρουβην, συμειων, λευι iudas: ifachar: zabulo: τ beniamin: da: et neph ioudas, ισαχαρ, αβουλων, και βενιαμιν, δαν, καινεφ θηλι: gad: τ aser. ioseph at erat in egypto. erat at θαλι, γαδ, και ασηρ. ιωσηφ δε ην εν αιγυπτω. ησασ δε οες sic q egressē fūt ex iacob: qnq τ sep τασαι ψυχαι αι εξελευσαι εε ιακωβ, τεντε και εβ tuaginta. mortuus eāt ioseph τ οες fra= do μηκοντα. ετελευτησε δε ιωσηφ και παντες οι αδελ tres eius: τ οis għatto illa. at filij israel cre= φοι αυτ, και πασα η γενεα εκεινη. οι δε υιοι ισραηλ ηυξη uerūt: τ multiplicati fūt: τ abidātes fuerūt: τ in= θισασ, και επι πληθυνθισασ, και χυδ αι οι εθνοντο, και κα ualuerūt valde multiplicati. multiplicauit at fra illos. sur τ ιχουον σφοδρα σφοδρα. επληθυνε δε η γη αυτ. ανε rexit at rex alter sup egyptu: q nō cognoscebat ση δε βασιλευς ετερος εω αυγυπτον, ος ουκ ηδει τον ioseph. dixit at genti sue: ecce gens fi iωσιφ. ειπαε δε τω εθνε αυτου, ιδου γο εθνος των υι lioz israel magna valde multitudo: τ pualet sup nos. ων ισραηλ μεγα πολυ πληθος, και ιχυει υπερ η= venite ergo sapiēter opprimatꝫ eos: ne forte μάς. δευτε ουν κατασοφισμεθα αυτους, μη ποτε multiplieēt: τ qñ acciderit nobis bellū: ad= πληθυνθη, και ηνικα αν συμβη ημιν πολεμος, προ= detur τ isti ad aduersarios. τ δε στεθησον εμ και αυτοι προς τους υπεναντιους. και εκ bellantes nos: egrediētur de fra. τ pre= πολεμισαστες ημας, εξελευσον εμ εκ της γης. και εωε fecit eis pfectos operū: vt affligerēt σησεν αυτοις επιστατας των εργαων, ινα κακωσωσιν eos in operibꝫ. τ edificaueēt ciuitates munitas αυτ. εν τοις εργαοις. και ακοδ ομοσασ πολεις θυρατας pharaōi: τ phihth et ramesse: τ on: q eit τω φαραω, την τε φιδωμ και ραμεση, και ων, η εστι η hellopolis. quāto at eos humiliabāt: tāto λιουπολις. και οτι δε αυτους εταπεινον, τοσούτω plures fiebant. τ inualuerūt valde. τ abhominatio πλειους εγινοντο. και ιχυον σφοδρα. και εβδελυσον ne huerūt egyptij a filijs israel. τ oppresse= το οι αιγυπτιοι ακωε των υιων ισραηλ. και κατεδυνα= rūt egyptij filios israel vi: τ afflix= σευον οι αιγυπτιοι τς υιους ισραηλ βια, και κατωδ υ= rūt eoz vitā in operibꝫ duris in των αυτων την ζωην εν τοις εργαοις τοις σκληροις, εν τω into τ lateritio: τ oibꝫ operibꝫ q in a= πικλω και τη laterithia, και πασι τοις εργαοις τοις εν τοις πε gris: fm oia opera: qbꝫ in fuit ut redegeēt eos διοις, και τα πάντα τα εργα, ων κατεδουλουντο αυτους

Transla. Chal. Incipit liber Exodus. Cap. i. Hec sunt nomina filioꝝ israel: qui ingressi sunt in egyptū cum iacob: singuli cum viris domus sue introierūt. Ruben: symeon: leui: et iudas: isachar: zabulon: et beniamin: dan: et nephtalim: et aser. Erant autem aie egresserunt de femore iacob septuaginta. Ioseph autem in egypto erat: et post mortem eius et omnibus fratribus eius: et cognatione sua: filii israel creuerunt: et quasi germinantes multiplicati sunt: et roborati sunt nimis: et impleverunt terram. Surrexit autem post haec rex novus super egyptum: qui nō cognoscebat ioseph. et ait ad populum suum: Ecce populus filiorum israel multus est et fortior nobis est. Venite ergo sapienter opprimamus eum: ne forte multiplicetur: et si in gruerit contra nos bellum: ad datur inimicis nostris: expugnatiſque nobis egre diatur de terra. Prepositi sunt itaque eis magistros operum: ut affligerent eos oneribus. Edificaverunt quoque: vrbes tabernaculoꝝ rum pharaōi phiton & ramesse. quātoque op primebāt eos: tāto magis multiplicabātur: & creſcebat. O derātq: filios israel egyptij: & affligēbant illudētes & inuidentes eis: atq ad amaritudinem perdūcebant vitam corū: operibus duris luti ooo & lateris: p omniq: fa= mulatu oooooooooooooo & lateris: p omniq: fa= mulatu oooooooooooooo quo in terre operibus ooooooooooooooooooooo

Polyglotte Bibeln

Polyglotte (mehrsprachige) Bibeln enthalten die biblischen Texte in Griechisch, Latein, Hebräisch und manchmal anderen Sprachen zum Zwecke der Textkritik.

Motivation heute

In der Computerlinguistik sind
parallele Korpora besonders
interessant für:

- mehrsprachige Lexikographie
- maschinelle Übersetzung

Parallele Korpora

Parallele Korpora unterscheiden sich durch:

- Anzahl der Sprachen (bi-/multilingual)
- Richtung der Alignierung
 - unidirektional (A \rightarrow B)
 - bidirektional (A \leftrightarrow B)
- Art der Alignierung
 - satzbasiert (sentence alignment)
 - wortbasiert (word alignment)

Alignierung

Ziel der Alignierung ist es die korrespondierenden Sätze oder Wörter der verschiedenen sprachigen Texte zu finden und zuzuordnen.

Dabei können mehrere Sätze/Wörter auch zu einem Satz/Wort zugeordnet werden.

“Gale & Church”-Algorithmus

*“A Program for Aligning Sentences
in Bilingual Corpora” (1991)*

William Gale und Kenneth Church
(beide AT&T Bell Laboratories)

“Gale & Church”-Algorithmus

Statistischer Ansatz mit der Grundidee, dass die Länge von korrespondierenden Sätzen korreliert.

Der korrespondierende Satz eines langen/kurzen Satzes in Sprache A ist in Sprache B ebenfalls lang/kurz.

Parallele Korpora: Satzebene

“Gale & Church”-Algorithmus

Die Wahrscheinlichkeit für eine Korrelation ergibt sich aus dem Verhältnis der Zeichenlängen der Sätze und deren Varianz.

Absätze müssen klar markiert sein und sich eindeutig entsprechen. Annahme von “Hard Boundaries” (Absätze) und von “Soft Boundaries” (Sätze).

“Clue Alignment”-Algorithmus

*“Combining Clues for Word
Alignment” (2003)*

Jörg Tiedemann

(Universität Uppsala, Schweden)

Parallele Korpora: Wortebene

“Clue Alignment”-Algorithmus

Ebenfalls ein statistischer Ansatz.
Es wird versucht sog. “Clues” zu finden um die Assoziierung von Wörtern in Quell- und Zielsprache zu bestimmen.

“Clues” können dabei z.B. Frequenz, Wortart, Phrasentypen oder die konkrete Wortform sein.

Parallele Korpora: Wortebene

“Clue Alignment”-Algorithmus

- haben die Wörter viele identische Zeichenfolgen?
- haben sie eine ähnliche relative Frequenz?
- treten sie zusammen auf (Kookurrenz)?
- wie groß ist der Positionsunterschied der Wörter?
- sind die Wortarten “kompatibel”?

Satzalignierte Korpora

European Parliament Proceeding
Parallel Corpus 1996-2003

- ca. 20 Millionen Wörter
- ca. 740.000 Sätze pro Sprache
- 11 Sprachen der EU

<http://www.isi.edu/~koehn/euoparl>

Satzalignierte Korpora

Aligned Hansards of the 36th
Parliament of Canada

- ca. 1,3 Millionen Wortpaare
- Englisch - Französisch

<http://www.isi.edu/natural-language/download/hansard/>

Satzalignierte Korpora

OPUS corpus

- übersetzte Web-Texte
- ca. 500.000 Wörter
- 5 Sprachen (FR, ES, SE, DE, JP)
- unidirektional

<http://logos.uio.no/opus/>

Satzalignierte Korpora

Slovene-English Parallel Corpus

- ca. 1 Millionen Wörter
- Englisch <-> Slowenisch
- bidirektional

<http://nl.ijs.si/elan/>

Wortalignierte Korpora

CRATER Multilingual Aligned Annotated Crp.

- drei Sprachen: EN, FR, ES
- 3 x 1 Millionen Token
- morphosyntaktisch annotiert, lemmatisiert
- unidirektional

www.comp.lancs.ac.uk/linguistics/crater/corpus.html

Parallele Korpora: Programme

Manatee/Bonito

Manatee (Server) verarbeitet nur vertikalen Text und Bonito (Client) zeigt Konkordanzen zu Querys an. Kann lediglich von bereits alignierten Sätzen Entsprechungen anzeigen.

Parallele Texte sind bei Manatee technisch verschiedene Korpora!

<http://www.textforge.cz/>

Parallele Korpora: Programme

Uplug

Uplug ist ein Korpus-Manager und kann konsolen- und web-basiert verwendet werden. Kann parallele Texte automatisch auf Satz- und Wortebene alignieren! Führt Preprocessing durch und bindet externe Tagger ein. Zeigt natürlich auch Konkordanzen zu Querys an.

<http://stp.ling.uu.se/cgi-bin/joerg/Uplug>