


NORDISKE STUDIER I LEKSIKOGRAFI

Titel:	SAOLHist – alla upplagor av SAOL i en och samma databas	
Forfatter:	Louise Holmer	
Kilde:	Nordiska Studier i Lexikografi 11, 2012, s. 287-295 Rapport från Konferens om lexicografi i Norden, Lund 24.-27. maj 2011	
URL:	http://ojs.statsbiblioteket.dk/index.php/nsil/issue/archive	

© Nordisk forening for lexicografi

Betingelser for brug af denne artikel

Denne artikel er omfattet af ophavsretsloven, og der må citeres fra den. Følgende betingelser skal dog være opfyldt:

- Citatet skal være i overensstemmelse med „god skik“
- Der må kun citeres „i det omfang, som betinges af formålet“
- Ophavsmanden til teksten skal krediteres, og kilden skal angives, jf. ovenstående bibliografiske oplysninger.

Søgbarhed

Artiklerne i de ældre Nordiske studier i lexicografi (1-5) er skannet og OCR-behandlet. OCR står for 'optical character recognition' og kan ved tegngenkendelse konvertere et billede til tekst. Dermed kan man søge i teksten. Imidlertid kan der opstå fejl i tegngenkendelsen, og når man søger på fx navne, skal man være forberedt på at søgningen ikke er 100 % pålidelig.

SAOLHist – alla upplagor av SAOL i en och samma databas

Louise Holmer

The Swedish Academy Glossary (SAOL), has so far been published in 13 editions since 1874. The goal of the SAOLHist project is to put all editions of SAOL, some of which exist only on paper, into a single database.

The number of lemmas of the Glossary has varied substantially over the years. The first edition comprised about 35,000 lemmas, while the latest has 125,000. In every new edition, new words have been added and obsolete ones have been excluded.

With SAOLHist, it will be possible to search for a word's timeline. Since all editions are searched simultaneously, information about the year of the first appearance of a word in the Glossary, or when it was subsequently removed, will be readily available.

Nyckelord: SAOL, SAOLHist, historisk lexikografi, svenskans ordförråd

1. Inledning

Sedan ett par år tillbaka pågår projektet SAOLHist vid redaktionen för *Svenska Akademiens ordlista* (SAOL) i Göteborg. Projektet går ut på att digitalisera alla de gamla upplagorna av ordlistan och sammanföra dem med de senare upplagorna i en databas som ska göras allmänt tillgänglig på nätet. De gamla upplagorna har skannats in och bearbetats digitalt vid redaktionen. Ansvarig för projektet är Sven-Göran Malmgren. Övriga redaktionsmedarbetare i SAOLHist är Daniel Berg, Monica von Martens och Louise Holmer. Skanning och bearbetning har utförts av Karin Malmgren, Joel Rogström och Oscar Sönnergren.

I den här artikeln presenteras projektet SAOLHist. I början beskrivs och exemplifieras det stora och diversifierade material som ligger till grund för databasen, och senare beskrivs digitaliseringsprocessen med den skanning och korrekturläsning som arbetet har inneburit. Avslutningsvis kommenteras projektet i sin helhet och tänkbara forskningsområden med anknytning till projektet behandlas.

2. Bakgrund till projektet SAOLHist

SAOL publicerades första gången 1874 och har sedan dess sammanlagt kommit ut i tretton upplagor, den senaste publicerad 2006. SAOL ger uppgift om uppslagsordens stavning och böjning och kommenterar i förekommande fall bruklighet och eventuellt det fackområde som uppslagsordet hör till. Däremot är det ingen definitionsordbok även om ungefär en femtedel av orden är försedda med definition eller någon annan sorts kommentar. Lemmaantalet är betydligt större än i andra enspråkiga svenska ordböcker, vilket framför allt beror på ett fylligt urval av sammansättningar (se vidare inledningskapitlet i SAOL13). För en fördjupning i bakgrunden till SAOL och ordlistans kännetecken hänvisas till exempelvis Gellerstam (red., 2009) och Johannisson & Mattsson (1974).

De olika upplagornas lemmaantal har varierat relativt mycket under åren, vilket åskådliggörs i tabell 1 nedan. Upplaga 2–5 utgör i stort sett bara nytryckningar av upplaga 1, medan upplaga 6 är den första som skiljer sig från föregångarna, framför allt genom sitt högre lemmaantal. Därför har vi låtit skanna in upplaga 1, 6, 7, 8, 9 och 10, medan upplaga 11–13 redan finns tillgängliga för redaktionen i digitalt format.

Upplaga	År	Huvudredaktör	Lemmaantal
1–5	1874–1883	F.A. Dahlgren	35 000
6	1889	F.A. Dahlgren	40 000
7	1900	Otto Hoppe	71 000
8	1923	Ebbe Tuneld	85 000
9	1950	Pelle Holm	155 000
10	1973	Gösta Mattsson	135 000
11	1986	Martin Gellerstam	115 000
12	1998	Martin Gellerstam	120 000
13	2006	Martin Gellerstam	125 000

Tabell 1. SAOL:s upplagor med årtal, huvudredaktör och lemmaantal (efter Gellerstam 2009: 56).

De olika upplagorna och deras respektive omfång och utformning har styrts såväl av redaktörerna och deras principer som av Svenska Akademiens önskemål. Antalet lemman är något som kommenteras i förordet till varje upplaga, och de flesta redaktörerna beskriver hur de har mönstrat ut föråldrade ord och tagit in nya. SAOL har hela tiden haft som uppgift att spegla det moderna ordförrådet. Det gör att en stor del av arbetet inför varje ny upplaga har bestått i att

lägga in nya ord såväl som att mönstra ut föråldrade ord (se vidare Gellerstam 2009; inledningen till SAOL13 2006; Berg, Holmer & Sköldberg 2010). Ordböcker i allmänhet tenderar att öka i omfång för varje ny upplaga, åtminstone tills något slags tak är nått beträffande sidantal, ordbokens vikt och liknande (Atkins & Rundell 2008: 22 f.). Detta stämmer till stor del för SAOL – ordlistans fysiska tjocklek har ökat med åren och lemmaantalet har ökat sedan den första upplagan. Den första upplagan var den lemmanumerärt minsta med ca 35 000 lemman, medan den nionde upplagan från 1950 är den hitintills mest omfattande med 155 000 lemman. I och med den elfte upplagan hamnade ordlistan i Göteborg med Martin Gellerstam som huvudredaktör, och där produceras den fortfarande. Arbetet med den fjortonde upplagan pågår nu under ledning av Sven-Göran Malmgren (huvudredaktör) och Sture Berg (biträdande redaktör).

3. Tidigare digitalisering av SAOL

SAOL har digitaliserats på olika sätt tidigare. Genom projekt Runebergs försorg finns upplaga 6 (1889) och 8 (1923) tillgängliga på internet, inskannade och delvis korrekturlästa (Projekt Runeberg 2007). Upplaga 11 (1986) finns i den allmänt tillgängliga Språkbankens konkordanser (Språkbanken 2011). Upplaga 12 (1998) gavs ut som CD-version och fanns också på internet under tidigt 2000-tal.

Upplaga 13 finns på CD och på nätet. CD-versionen heter SAOL Plus och kom 2007. Den har betydligt mer information och fler sökmöjligheter än pappersversionen. Tillägget ”Plus” betecknar framför allt att den innehåller samtliga *teoretiskt* tänkbara böjningsformer för varje uppslagsord, i motsats till den tryckta ordlistan som bara ger ett fåtal böjningsformer för de flesta ord och ofta inga alls vid sammansatta ord. Användaren kan i SAOL Plus välja att söka bland uppslagsord, i artikeltexten eller bland böjningsformer (se vidare Berg, Holmer & Hult 2008; Martola 2008).

Nätversionen av den trettonde upplagan publicerades i mars 2009 i en faksimilversion (Svenska Akademien 2011). Eftersom ordlistan bara ligger som bilder i pdf-format och inte i en databas går det inte att göra några förfinade sökningar såsom att söka på delar av ett ord eller i artikeltext eller liknande.

Under 2011 påbörjades också, på initiativ av Svenska Akademien, arbetet med att ta fram en app-version av SAOL13 för smarta mobiltelefoner. Appen är i skrivande stund inte färdig, men dess lansering är beräknad till november 2011.

4. Lemmaantalet och omsättningen av lemman mellan upplagorna

Överlappningen av lemman mellan de olika upplagorna är stor, men uppskattningsvis har ca 200 000 lemman vid ett eller flera tillfällen funnits med i ordlistan. Med projektet SAOLHist blir det möjligt att se i vilka upplagor ett visst ord har funnits med och därigenom när det kom in i ordlistan och när det (eventuellt) ströks ur ordlistan. Det blir också möjligt att se vilka ord som bara funnits med i en enda upplaga. Vidare kan man också få en uppfattning om ett ”basordförråd”, alltså lemman som funnits med i alla upplagor, från den första till den sista. Det handlar alltså om sådana ord som hittills inte har mönstrats ut av bl.a. åldersskäl, utan som har behållit sin aktualitet i språket och därmed i ordlistan.

För att få en fingervisning om hur lemmauppsättningen har ändrats mellan upplagorna gjordes en mycket liten undersökning vid redaktionen hösten 2009. Lemman på *Fa*-undersöktes i den första, nionde och trettonde upplagan. Totalt handlade det om 440 lemman, där 80 av dem fanns med i upplaga 1, 320 i upplaga 9 och 305 i upplaga 13. Av dessa var 50 lemman gemensamma för de tre upplagorna, t.ex. *fabel*, *fabrik*, *falsk* och *fagott*. Så många som 150 lemman av de 440 fanns inte längre kvar i SAOL13. Exempel på sådana utmönstrade ord är *fadersinstinkt*, *fajanskopp* och *fallvirke*.

Här bör tilläggas att varje enkelt (osammansatt) ord för med sig sammansättningar och avledningar. Ett exempel är lemmat *fabrikation* vilket saknas i första upplagan. I nionde upplagan finns *fabrikation* med och har ett flertal sammansättningar, t.ex. *fabrikationsfel* och *fabrikationshemlighet*. Varje nytt lemma för alltså också ofta med sig flera andra lemman.

Även om undersökningen som beskrivs ovan är minimal säger den ändå något om den stora omsättningen av lemman som har ägt rum mellan de olika upplagorna. Se tabell 2 för en enkel uppställning över lemmaförekomsten.

5. Utvecklingen av ordlistans artikelstruktur

Till och med 11:e upplagan placerades de flesta sammansättningar in under sina respektive huvudord, vilket gjorde att ord som etymologiskt hörde ihop också återfanns på samma ställe i ordlistan. Med 12:e upplagan infördes i stället en sorteringsordning som innebär att alla lemman placeras i bokstavsordning oavsett etymologi (Gellerstam 2009: 71–73, SAOL12: VI f.). Dessutom infördes den s.k. lemma-lexemmodellen, vilken innebär att tidigare olika lemman med samma böjning sorteras in under samma huvudord och blir olika lexem. Tidi-

gare hade homografer dessutom ofta sorterats in under samma lemma även om böjningen skilde dem åt (t.ex. *ljus* a. och *ljus* s., se SAOL12: X). I och med 12:e upplagan gjordes också en noggrannare genomgång av varje lemmas ordklass-tillhörighet (SAOL12: XIX).

Lemma	SAOL1	SAOL9	SAOL13
falsifierbar			X
falsifiering			X
falsifikat		X	X
falsifikation			X
falsk	X	X	X
falskant		X	
falskdeklarant		X	X
falskdeklaration		X	X
falskeligen	X	X	X
falskhet	X	X	X
falsklarm			X
falsklegg			X
falskmyntare		X	X
falskmyntarlīga		X	

Tabell 2. Exempel på lemmaförekomst vid *Fa-* i tre olika upplagor av SAOL.

Många av upplagorna uppvisar stora skillnader i sättet att presentera artiklarna. I de tidiga upplagorna anges exempelvis genus för substantiven. Ett ord som *björk* ser ut enligt följande i första respektive sjunde upplagan:

SAOL1 Björk (pl. *-ar*) s. f.
SAOL7 björk (*-en, -ar*) rf.

I första upplagan ges alltså först pluraländelsen och därefter ordklassmarkeringen ”s. f.” vilket står för substantiv, femininum. I sjunde upplagan ges förutom pluraländelsen även ändelsen i bestämd form. Ordklassangivelsen är också modifierad till ”rf.”, vilket står för reale, femininum (se vidare Malmgren 2002: 10–12). Inledningen till sjunde upplagan bjuder på följande förklaring:

Utanför de levande varelsernas område, t. e. i fråga om *mur, vägg*, nyttjar man i de bildades språk numera icke ofta *han* och *hon*; i stället säges *den*. Med andra ord: maskulinum och femininum hafva här blifvit ersatta med ett nytt genus, som på sista tiden fått namnet *realgenus*, och som nu i Ordlistan betecknas med *r*. (SAOL7: VIII)

Hur verben presenteras skiljer sig åt mellan de olika upplagorna. Faktum är att verb är den ordklass vars presentationsstruktur uppvisar störst och flest skillna-

der. För att åskådliggöra några av olikheterna mellan upplagorna har jag valt ett deponentiellt verb från första konjugationen, *hoppas*. Överst i uppställningen visas den senaste upplagan och sedan följer de andra i omvänd kronologisk ordning.

- 13, 12 **hoppas** v. *hoppades* till ¹*hopp1*
- 11 **hoppas** -ades v. dep.
- 10 **hopp|as** -ades v. dep.
- 9 **hopp|as** -ades itr. dep.
- 8 **hopp|as** (-ades) i. o. t. dep.
- 7 **hopp|as** (-ades) i. dep.
- 6 **Hoppas** (-ades) dep.
- 1 **Hoppas** (-as, -ades, -ats), v. dep.

Ett antal detaljer syns tydligt i uppställningen ovan. För det första är det lemmats utseende och presentation: i upplaga 1–6 användes versal initialt, något som togs bort i 7:an. I upplaga 7 infördes däremot lodstreck efter ordstammen, en detalj som fick vara kvar i flera upplagor. Från upplaga 11 och framåt finns inget lodstreck i den här typen av enkelt (icke sammansatt) uppslagsord.

I den tolfte och trettonde upplagan har varje verb, även om det är sammansatt, ordklassbeteckningen "v." Först ges ordklassuppgift i förkortad form och därefter preteritumformen *hoppades* fullt utskrivet. I alla de tidigare upplagorna ges bara böjningsändelsen i preteritum, eller som i första upplagans fall, böjningsändelserna i presens, preteritum och supinum (-as, -ades, -ats).

Den första och de sista fyra upplagorna har ordklassbeteckningen "v." för verb medan upplagorna däremellan har beteckningen "dep." för deponens och dessutom en transitivetsbeteckning: "i." för intransitivt och "t." för transitivt verb. Där förutsätts alltså läsaren veta vad förkortningarna står för och dessutom vad de innebär. I dagens ordlista är det oklart om man skulle nå förståelse hos användarna om man skrev "i. o. t. dep." i stället för "v.", men äldre tiders användare hade troligtvis lite bättre kunskap om grammatisk terminologi.

Definition saknas i alla upplagor. Upplaga 12 och 13 har däremot en hänvisning till substantivet *hopp* och dess första lexem ("1. förhoppning"), vilket indirekt antyder betydelsen av verbet *hoppas*.

6. Digitaliseringsarbetet

Skanningen och bearbetningen av det inskannade materialet har gjorts vid ordlisteredaktionen. Lite förenklat kan processen sägas ha gått till enligt följande: originalen har skannats in, tolkats av programmet Omnipage och därefter kontrollerats manuellt. Sedan har filerna exporterats till xml-format. Därefter går det

att se materialet i en databas eller visa dem på en webbsida. I ett inledande skede har endast lemmat och ordklassuppgifterna extraherats (i den mån ordklass finns med, vilket inte alltid är fallet i de äldre upplagorna).

Originalsidorna bjuder på en typografisk mångfald. Typsnitten varierar mellan upplagorna, bruket av kursiv och rak stil växlar liksom teckengraden inne i artikeltexten. Detta har medfört vissa svårigheter i efterarbetet. Tolkningen av texten i Omnipage har ibland haft svårigheter med att identifiera just det som det är programmerat att identifiera, nämligen ord med fetstil (lemmat och dess sammansättningar och avledningar) och ordklassbeteckning (som, vilket visats tidigare, varierar en hel del mellan upplagorna). Om då programmet ”missat” att hämta upp ett fetstilt ord som dessutom har många sammansättningar blir resultatet att efterleden kopplas ihop med helt fel förled, alternativt saknas helt. Se figur 1 för exempel på en sida ur upplaga 7 som skannats in, tolkats av programmet och korrekturlästs.

F ¶	
<p>f (<i>fet</i>; pl. =, best. <i>f-en</i>) n. <i>F.dur. f-dursskala. f-ljud. ¶</i></p> <p>fabel (-n; <i>fabler</i>) r. - aktig (~t; ~are) a. - bok. - diktare. - hjälte. - samling. - fabricer a (-ade) t. - ing rf. fabrik (-en; -er) m. fabriks alster. - anläggning. - arbetare. - arbete. - arbeterska. - byggnad. - distrikt. - drift. - flicka. - idkare. - inspek - tion. - inspektör. - märke. - mässig (-t; ~are) a. - ort. - pris. - rörelse. - skorsten. - stad. - vara. - fabrik ant (-en; -er) m. - at (-et; pl. = 1. ~er) n. - ation (-en; ~er) r. - ör (-en; ~er) m. ¶</p> <p>fabulös (-t; ~are) a. ¶</p> <p>facit n. - bok. - tabell. ¶</p> <p>fack (-et; pl. = n. - bildning. - förening. - in - sikt (er). - kunskap. - lärare. - man. - mäs - sig (~t; ~are) a. - skola. - skrift. - studium. - term. - tidskrift. - verk. <i>Korsvirke. ¶</i></p> <p>fackl a (-an; -or) rf. fackell belysning.</p>	<p>fadder (-n; <i>faddrar</i>) mf. - gåfva. - kyss. - fadderskap (-et; pl. =) n. ¶</p> <p>faddhet (-en; -er) rf. ¶</p> <p>fader, smdr. far (<i>fadem</i>, äfv. <i>fordren</i>; <i>fäder</i>, best. <i>fäderna</i>, äfv. <i>poet. fädren</i>) <i>m.</i> - lös. - mord. - mördare. - mörderska. ¶</p> <p>faders glädje. - hem <i>relig.</i> - hjärta. - hus <i>relig.</i> - kärlek. - namn. - sinne. - välde. - far bro de r. - fa de r. - föräldrar. - mo de r. farsgubbe (n). <i>skämts. o. hvard. - fader lig</i> (-t; ~are) a. - ligen (o. - ligt) adv. - skap (-et) n. - fadervår (pl. =) n. ¶</p> <p>fager (-t; <i>fagra; fagrare</i>) a. - hyllt (n. ¶ =) p. a. - kindad (-t) p. a. <i>poet.</i> ¶</p> <p>faggorna best. pl. <i>Hafva ngt i f.</i> ¶</p> <p>fagott (-en; -er) r. - blåsare (-n; pl. =) m. - stämma. fagottist (-en; -er) m. - fajans (-en; -er) r. - fabrik. - fat. - kärl. - fakir (-en; -er) m. ¶</p> <p>faksimile (-t; pl. = 1. -n) n. - maner. -</p>

Figur 1: Sida ur SAOL7 efter skanning, tolkning och kontrolläsning.

Markeringarna i figuren visar olika ställen där programmet har haft svårt att tolka tecknen, exempelvis vid lodstrecken i **fabriks**|**alster** och **faders**|**glädje** och vid tecknet ~ på ett flertal ställen, t.ex. efter **fabelaktig** och **fabriksmässig**.

I den färdiga databasen kommer man att kunna söka med modern stavning och med originalstavning. Det innebär att om man söker på lemmat *avhandla* ska man få resultatet att lemmat funnits med i alla upplagor, trots den tidigare stavningen *afhandla*. Vid lemmen med variantstavningar som t.ex. *schyst* och *juste* ska man kunna söka på båda stavningarna. I resultatlistan visas sedan båda stavningsvarianterna med respektive upplaga (*juste* kom in i den 8:e upplagan och stavningen *schyst* i den 11:e).

7. Vad betyder ordlistans variation för SAOLHist?

I artikeln har jag gett ett flertal exempel på hur varierad ordlistan har varit i sitt sätt att presentera artiklarna i de olika upplagorna. Vad får det för inverkan på projektet SAOLHist? För det första är det ett stort material att arbeta med för alla inblandade. Det gör att det praktiska arbetet såsom inskanning och efterarbete har varit relativt omfattande. Att sedan läsa slutkorrektur blir också ett mycket stort arbete även om delar av det kan automatiseras (t.ex. genom att de olika upplagorna körs mot varandra). De olika sätten att markera ordklass påverkar också. Många lemmen i de tidigare upplagorna saknar helt ordklassangivelse och många av dem som har ordklassmarkering har olika sätt att markera det på. Det leder till en sorts normalisering beträffande ett flertal av lemmanas ordklassuppgifter. Lemman som tidigare har haft ordklassmarkeringen ”i. o. t. dep.” för intransitivt och transitivt deponensverb får i stället bara ett ”v.” för verb.

8. Sammanfattning och framtida användningsområden för SAOLHist

I den här artikeln har jag beskrivit de olika upplagorna av SAOL och hur dessa har insortrats i projektet SAOLHist. I och med SAOLHist kommer i princip samtliga upplagor att finnas sökbara i en och samma databas. Även om SAOL har givits ut i elektronisk form tidigare blir det här det första projekt som samlar alla upplagor på ett ställe. Med SAOLHist kommer man att kunna söka fram ord som funnits med i alla upplagor såväl som ord som bara funnits med i en upplaga innan de har blivit utmönstrade. Totalt handlar det om ca 200 000 olika lemmen, däribland ett basordförråd som utgörs av ord som funnits med i alla upplagor.

Förhoppningen är att databasen ska kunna ligga till grund för både språk-historisk, lexikografisk och lexikologisk forskning och att även allmänheten kommer att kunna ha glädje av den. Kanske blir också en del ordlistanvändare

mindre upprörda när ord stryks ur den aktuella upplagan av ordlistan, eftersom SAOLHist kommer att utgöra en permanent viloplats för utmönstrade ord.

KÄLLOR OCH LITTERATUR

- Atkins, B.T. Sue & Michael Rundell, 2008: *The Oxford Guide to Practical Lexicography*. Oxford.
- Berg, Sture, Louise Holmer & Ann-Kristin Hult, 2008: SAOL Plus – a New Swedish Electronic Dictionary. I: Bernal, E. & J. DeCesaris (eds.), *Proceedings of the XIII Euralex International Congress (Barcelona 15–19 July 2008)*, S. 291–296. (Cd-rom.)
- Berg, Sture, Louise Holmer & Emma Sköldberg, 2010: Time to say goodbye? On the exclusion of solid compounds from the Swedish Academy Glossary. I: Dykstra, Anne & Tanneke Schoonheim (eds.), *Proceedings of the XIV Euralex International Congress (Leeuwarden, 6–10 July 2010)*, S. 567–576
- Gellerstam, Martin, 2009: SAOL i många upplagor. I: Gellerstam, Martin (red.), S. 53–83.
- Gellerstam, Martin (red.), 2009: SAOL och tidens flykt. Några nedslag i ordlistans historia. Stockholm.
- Johannisson, Ture & Gösta Mattsson, 1974: Svenska Akademiens ordlista under 100 år. Skrifter utgivna av Svenska språknämnden 55. Stockholm.
- Malmgren, Sven-Göran, 2002: Normering i Svenska Akademiens ordlista 1874–1950: principer och resultat. I: *LexicoNordica* 9. S. 5–20.
- Martola, Nina, 2008: SAOL Plus – SAOL på cd-rom. I: *LexicoNordica* 15. S. 261–278.
- Projekt Runeberg 2007: Svenska Akademiens Ordlista. <http://runeberg.org/saol/> (september 2011)
- SAOL1–13 = Svenska Akademiens ordlista 1–13 uppl., 1874–2006.
- SAOL Plus = Svenska Akademiens ordlista över svenska språket. 13 uppl. Cd-rom. 2007: Stockholm.
- Språkbanken 2011: Konkordanser – SAOL 11. <http://spraakbanken.gu.se/konk/> (september 2011)
- Svenska Akademiens 2011: Svenska Akademiens ordlista. <http://www.svenskaakademien.se/>
- svenska_spraknet/svenska_akademiens_ordlista/saol_pa_natet/ordlista (september 2011)

Louise Holmer

Lexikaliska institutet, Inst. för svenska språket, Göteborgs universitet
 louise.holmer@svenska.gu.se