



Identifying Literary Texts with Bigrams

ANDREAS.VAN.CRANENBURGH@Huygens.knaw.nl
CORINA KOOLEN — C.W.Koolen@UvA.nl



SURVEY

As part of the project *The Riddle of Literary Quality*, a large (13k respondents) survey was held of the Dutch reading public (Cf. www.hetnationalelezeronderzoek.nl). Participants rated books they had read on two scales:

LITERARINESS 1–7 scale
BAD/GOOD 1–7 scale

METHOD

- linear Support Vector Machine (SVM)
- trained on bigram counts and mean of ratings.
- 10-fold cross-validation
- Tasks:

CLASSIFICATION binary classification after applying threshold to scale (predict high vs. low rating)
REGRESSION predict numeric rating

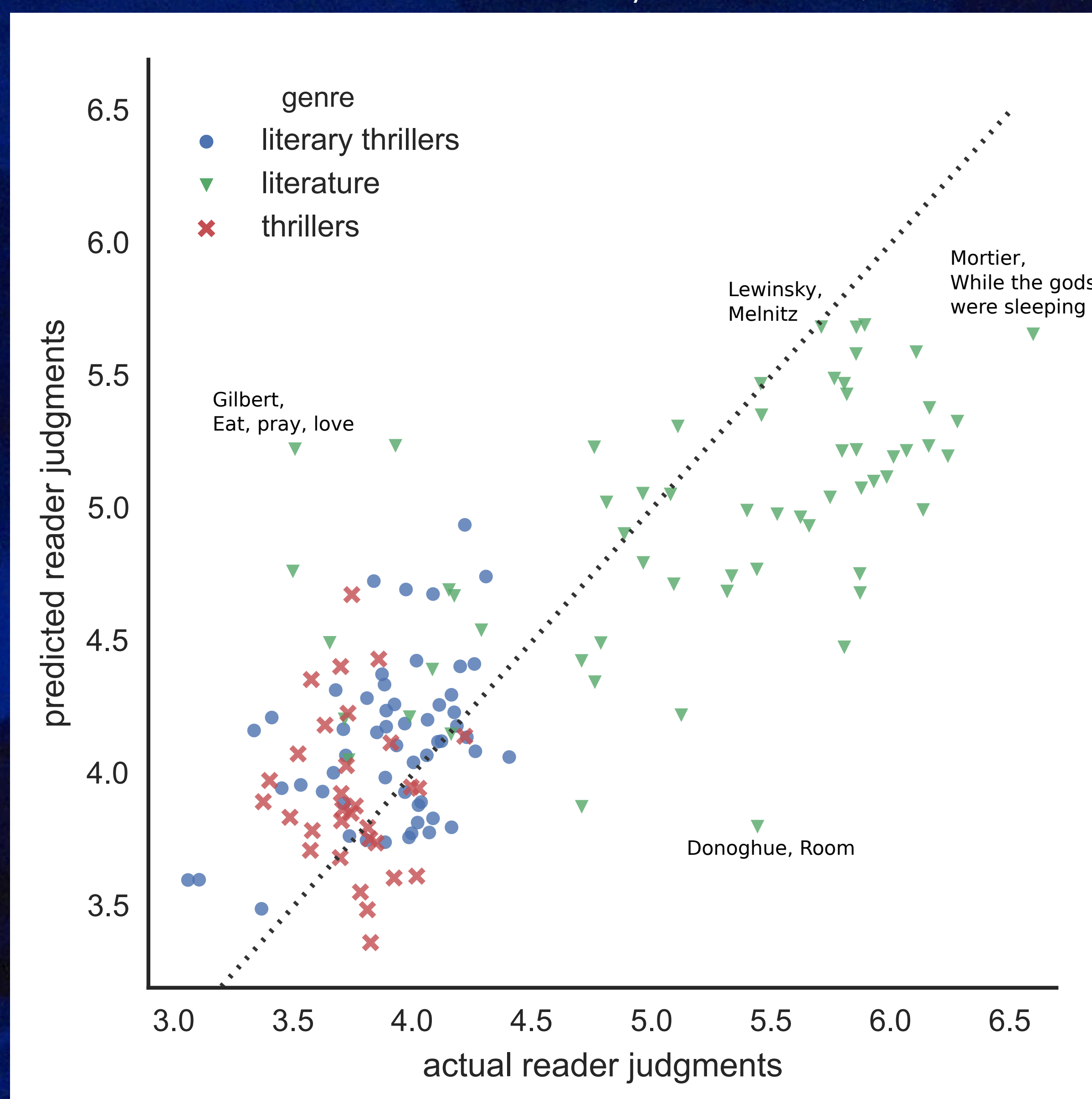
RESULTS

Classification results (accuracy; percentage correct):

Features	Literary	Bad/good
Content	90.4	63.7
Style	89.0	63.0

Evaluation of the regression models; R^2 scores (percentage of variation explained), root mean squared error in parentheses (1–7).

Features	Literary	Bad/Good
Content	61.3 (0.65)	33.5 (0.49)
Style	57.0 (0.67)	22.2 (0.52)



Regression results for literariness using content bigrams.

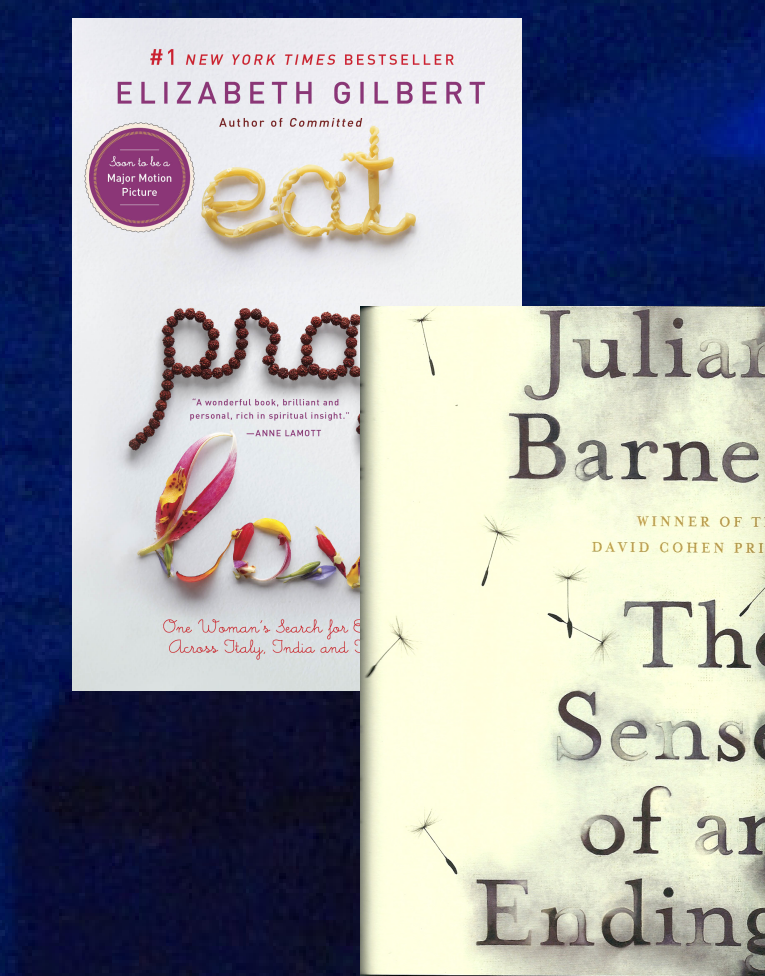
TEXTS

A subset (cf. Corpus box) of the Riddle corpus of 400 recent, successful novels, containing thrillers, literary thrillers, and literary novels.

Features: Bag-of-Words model based on bigrams that appear in 60–90 % of texts.

CONTENT BIGRAMS all words.

STYLE BIGRAMS function words plus content words replaced with POS tags.



CONCLUSIONS

Perceived literariness can be predicted, to a surprisingly large extent.

Literariness is easier to predict than general quality.

Both content features and style features perform well.

The regression results show that the model does not just recognize genres, but also accounts for literariness differences *within* genres.

How to describe the nature of literariness still an open question.

INTERPRETATION

Weights of the top 20 most important content features (above) and top 5 most important style features (below), indicative of literary (left), and non-literary texts (right).

weight	literary features, content	weight	non-literary features, content		
12.1	<i>de oorlog</i>	the war	-6.1	<i>de moeder</i>	the mother
8.1	<i>het bos</i>	the forest	-5.1	<i>keek op</i>	looked up
8.1	<i>de winter</i>	the winter	-4.9	<i>mijn hoofd</i>	my head
6.6	<i>de dokter</i>	the doctor	-4.9	<i>haar moeder</i>	her mother
5.8	<i>zo veel</i>	so much	-4.7	<i>mijn ogen</i>	my eyes
4.8	<i>nog altijd</i>	yet still	-4.7	<i>ze keek</i>	she looked
4.5	<i>de meisjes</i>	the girls	-4.5	<i>mobilele telefoon</i>	mobile telephone
4.3	<i>zijn vader</i>	his father	-4.2	<i>de moord</i>	the murder
4.0	<i>mijn dochter</i>	my daughter	-4.0	<i>even later</i>	a while later
3.9	<i>het boek</i>	the book	-3.8	<i>nu toe</i>	(until) now
3.8	<i>de trein</i>	the train	-3.5	<i>zag ze</i>	she saw
3.7	<i>hij hem</i>	he him	-3.4	<i>ik voel</i>	I feel
3.7	<i>naar mij</i>	at me	-3.3	<i>mijn man</i>	my husband
3.5	<i>zegt dat</i>	says that	-3.2	<i>tot haar</i>	to her
3.5	<i>het land</i>	the land	-3.2	<i>het gebouw</i>	the building
3.5	<i>een sigaret</i>	a cigarette	-3.2	<i>liep naar</i>	walked to
3.4	<i>haar vader</i>	her father	-3.1	<i>we weten</i>	we know
3.4	<i>een boek</i>	a book	-3.1	<i>enige wat</i>	only thing
3.2	<i>de winkel</i>	the shop	-3.1	<i>en dus</i>	and so
3.1	<i>elke keer</i>	each time	-3.0	<i>in godsnaam</i>	in god's name

weight	literary features, style	weight	non-literary features, style		
21.8	<i>! WW</i>	! VERB ,	-13.8	<i>nu toe</i>	until now
20.5	<i>u ,</i>	you (FORMAL) ,	-13.4	<i>en dus</i>	and so
18.0	<i>haar haar</i>	her (her)	-13.4	<i>achter me</i>	behind me
16.5	<i>SPEC :</i>	NAME :	-13.2	<i>terwijl ik</i>	while I
15.4	<i>worden ik</i>	become I	-13.1	<i>tot nu</i>	until now

- Some features relate to topics of genres (*war, murder*)
- 'Timeless' character of literary features (*book, letter*) vs. modern technology in non-literary books (*computer, mobile phone*).
- More markers of colloquial language in non-literary texts, in both the content and style features

REFERENCES

This work is part of The Riddle of Literary Quality, a project supported by the Royal Netherlands Academy of Arts and Sciences through the Computational Humanities Program.

- PAPER Andreas van Cranenburgh, Corina Koolen (2015). Identifying Literary Novels with Bigrams. Workshop on Computational Linguistics for Literature, Denver, Colorado, USA.
- PROJECT The Riddle of Literary Quality; cf. <http://literaryquality.huygens.knaw.nl>
- PAINTING Christine Bittremieux (2007). Untitled. Oil on canvas. 60 x 80 cm. www.bittremieux.nl

CORPUS

Legend: literary novel, thriller, literary thriller.

- APPEL, RENÉ *Van twee kanten*
- AUEL, JEAN MARIE *Het lied van de grotten*
- AVALLONE, SILVIA *Staal*
- BALDACCIO, DAVID *De provocatie, De rechtvaardigen, De zesde man, Die zomer, Familieverraad, Geniaal geheim, In het geheim, Onschuldig, Rechtfeloos, Verloos ons van het kwaad*
- BERNLEF, J. *De een zijn dood, Geleende levens*
- BINET, LAURENT *Hinh*
- BLUM, JENNA *In tveestijd*
- CLANCY, TOM *In het vizier, Op leven en dood, De ogen van de vijand*
- COBEN, HARLAN *Lavenslijn*
- CRONIN, JUSTIN *De oversteek*
- CUSSLER, CLIVE *Dodenschap, Wassende maan*
- DUIKZEUL, LIENEKE *Koudse lente, Verloren zoon*
- Die, ADRIAAN VAN *Tikkop*
- DONOGHUE, EMMA *Kamer*
- DORRSTEIN, RENATE *De stiefmoeder*
- DURLACHER, JESSICA *De held*
- ENQUIST, ANNA *De verdoovers*
- EVANS, NICHOLAS *De vergeving*
- FORBES, ELENA *Sterf met mij*
- FORSYTH, FREDERICK *De cobra*
- FRAGOSO, MARGAUX *Tijger, Ijger*
- FRENCH, NICCI *Blauwe maandag, Medeplichtig*
- GALEN, ALEX VAN *Suskind*
- GEORGE, ELIZABETH *Een duister vermoeden*
- GERRITSEN, TESS *De Mefisto Club, Het aandenken, Het stille meisje, Koud hart, Sneeuwval, Verdwin*
- GILBERT, ELIZABETH *Eten, bidden, beminnen*
- GIORDANO, PAOLO *De eenzaamheid van de priemgetallen*
- GIPHART, RONALD *Usland*
- GUDENKAUF, HEATHER *In stilte gehuld*
- HANNAH, SOPHIE *Kleine meid, Moederziel*
- HARBACH, CHAD *De kunst van het veldspel*
- HART, MAARTEN 't *Verlovingstijd*
- HAYNES, ELIZABETH *Waarheen je ook vlucht*
- HEIJDEN, A.F.Th. VAN DER *Tonio*
- HILL, LAWRENCE *Het negerboek*
- HOAG, TAMI *Dieper dan de doden*
- HODGKINSON, AMANDA *Britannia Road 22*
- HOSSEINI, KHALED *Dulzend schitterende zonnen*
- IRVING, JOHN *De laatste nacht in Twisted River, In een mens*
- JAMES, ERICA *Schaduwleven, Zussen voor altijd*
- JANSSEN, ROEL *De tiende vrouw*
- JAPIN, ARTHUR *De overgave, Vaslav*
- KEPLER, LARS *Hypnose*
- KLUUN *Haantjes*
- KOCH, HERMAN *Het diner, Zomerhuis met zwembad*
- KORYTA, MICHAEL *Begraven*
- KRAUSS, NICOLE *Het grote huis*
- KROONENBERG, YVONNE *De familie light blues*
- LACKBERG, CAMILLA *Ijsprinses*
- LAPIDUS, JENS *Bloedlink, Snel geld, Val dood*
- LARSSON, STIEG *De vrouw die met vuur speelde, Gerechtigheid, Mannen die vrouwen haten*
- LAUNSPACH, RIK *1953 (De Storm)*
- LEWINSKY, CHARLES *De verborgen geschiedenis van Courtillon, Het lot van de familie Meijer*
- LÄCKBERG, CAMILLA *Steenhouwer, Vuurforenwachter*
- MASTRAS, GEORGE *Iranen over Kashmir*
- MCCOY, SARAH *De bakkersdochter*
- McFADYEN, CODY *Tijd om te sterven*
- McNAB, ANDY *Oorlogswond*
- MOOR, MARENTE DE *De Nederlandse maagd*
- MOOR, MARGRIET DE *De schilder en het meisje*
- MORTIER, ERWIN *Godenslaap*
- NESBO, Jo *De schim, De sneeuwman, Het pantserhart*
- NOORT, SASKIA *Afgunst, De eetclub, De verbouwing, Koorts*
- PATTERSON, JAMES *De affaire, Hifite*
- PAUW, MARION *Daglicht, Jetset*
- PICK, ALISON *Donderdagkind*
- PICOUIT, JODI *Negenfien minuten*
- ROBERTS, NORA *Erbetoan*
- ROBOTHAM, MICHAEL *Gebraken*
- ROSE, KAREN *Moord voor mij*
- ROSENBOOM, THOMAS *Zoete mond*
- ROSENFELDT, HJORTH *Wat verborgen is*
- ROSNAY, TATIANA DE *Het appartement, Het huis waar jij van hield*
- RUIZ ZAFÓN, CARLOS *De gevangene van de hemel, De schaduw van de wind*
- SANSOM, CHRISTOPHER JOHN *Winter in Madrid*
- SCHOLTEN, JAAP *Kameraad Baron*
- SIEBELINK, JAN *Het lichaam van Clara, Oscar*
- SLAUGHTER, KARIN *Genadeloos, Genesis, Gevallen, Onaantastbaar, Verbroken*
- STEVENS, CHEVY *Vermist*
- TERLOUW, JAN / TERLOUW, SANNE *Hellehanden*
- TRUSSONI, DANIELLE *Het uur van de engelen*
- VERHOEF, ESTHER *Alles te verliezen, Close-up, Déja vu, Erken mij, Tegenlicht*
- VERHULST, DIMITRI *De laatste liefde van mijn moeder*
- VERMEER, SUZANNE *Après-ski, Cruise, De suite, De vlucht, Zomertijd*
- VISSER, JUDITH *Stuk*
- VLUGT, SIMONE VAN DER *Blauw water, Het laatste offer, In mijn dromen, Op klaarlichte dag*
- WATSON, S.J. *Voor ik ga slapen*
- WINTER, LEON DE *Recht op terugkeer, VSV of daden van onbaatzuchtigheid*
- YALOM, IRVIN D. *Het raadsel spinoza*