

Computational Argumentation Quality Assessment in Natural Language

Henning Wachsmuth Bauhaus-Universität Weimar

Nona Naderi University of Toronto

Yonatan Bilu IBM Research – Haifa

Yufang Hou IBM Research – Ireland

Vinodkumar Prabhakaran Stanford University

Tim Alberdingk Thijm University of Toronto

Graeme Hirst University of Toronto

Benno Stein Bauhaus-Universität Weimar

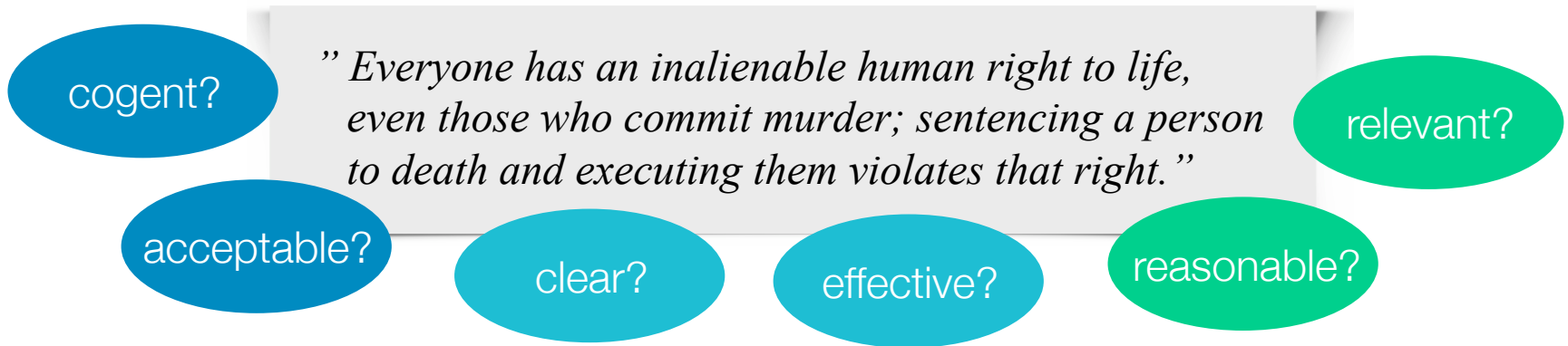
**Bauhaus-Universität
Weimar**

www.webis.de

henning.wachsmuth@uni-weimar.de

Motivation

- **Argument mining**
 - Identifies arguments in natural language text
 - Does not assess quality
- **Argumentation quality assessment**
 - Critical for any application built upon argument mining



- **Challenges**
 - Several quality dimensions at different granularities
 - Some highly subjective
 - How *should* we argue vs. how *do* we argue

Background

▪ Debating Technologies at Dagstuhl

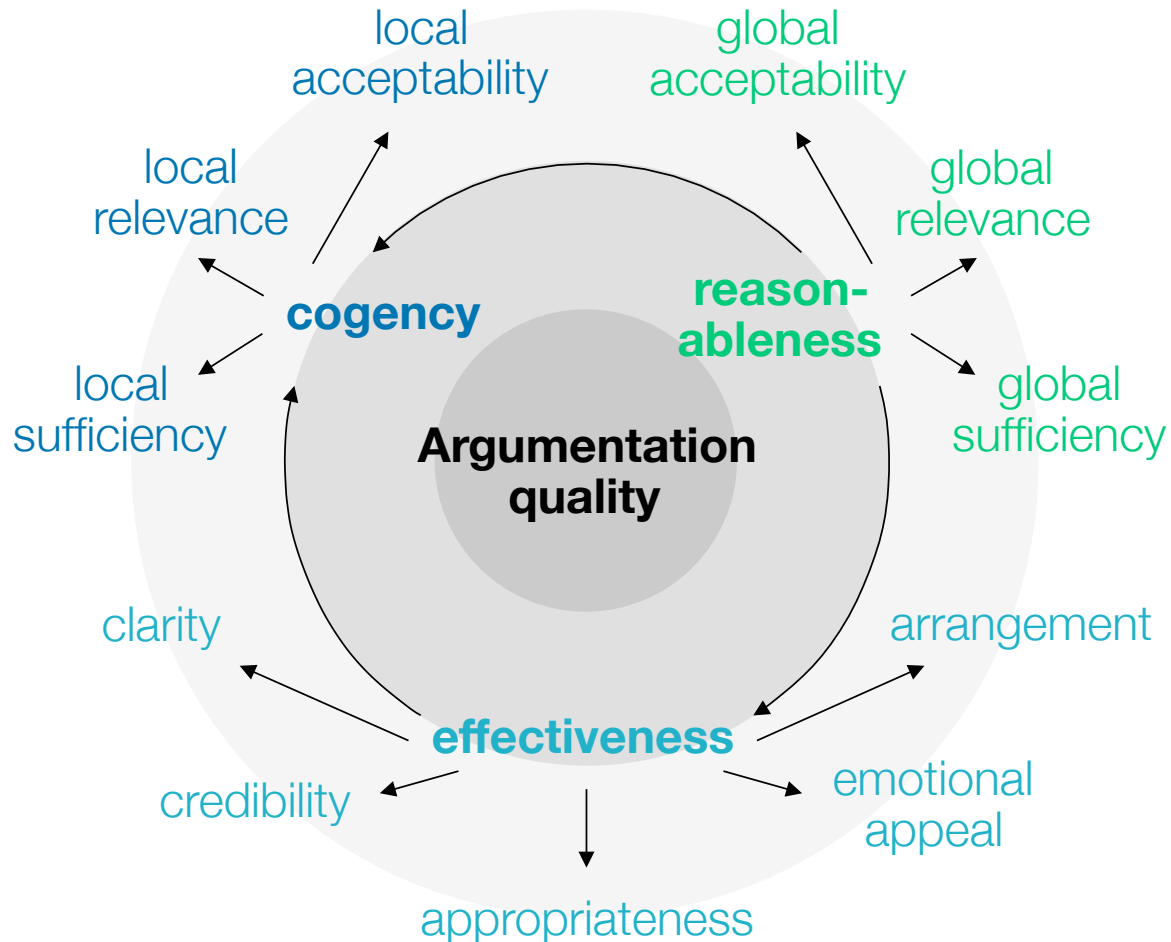
- Common understanding of argumentation quality missing
- Working group to coordinate research



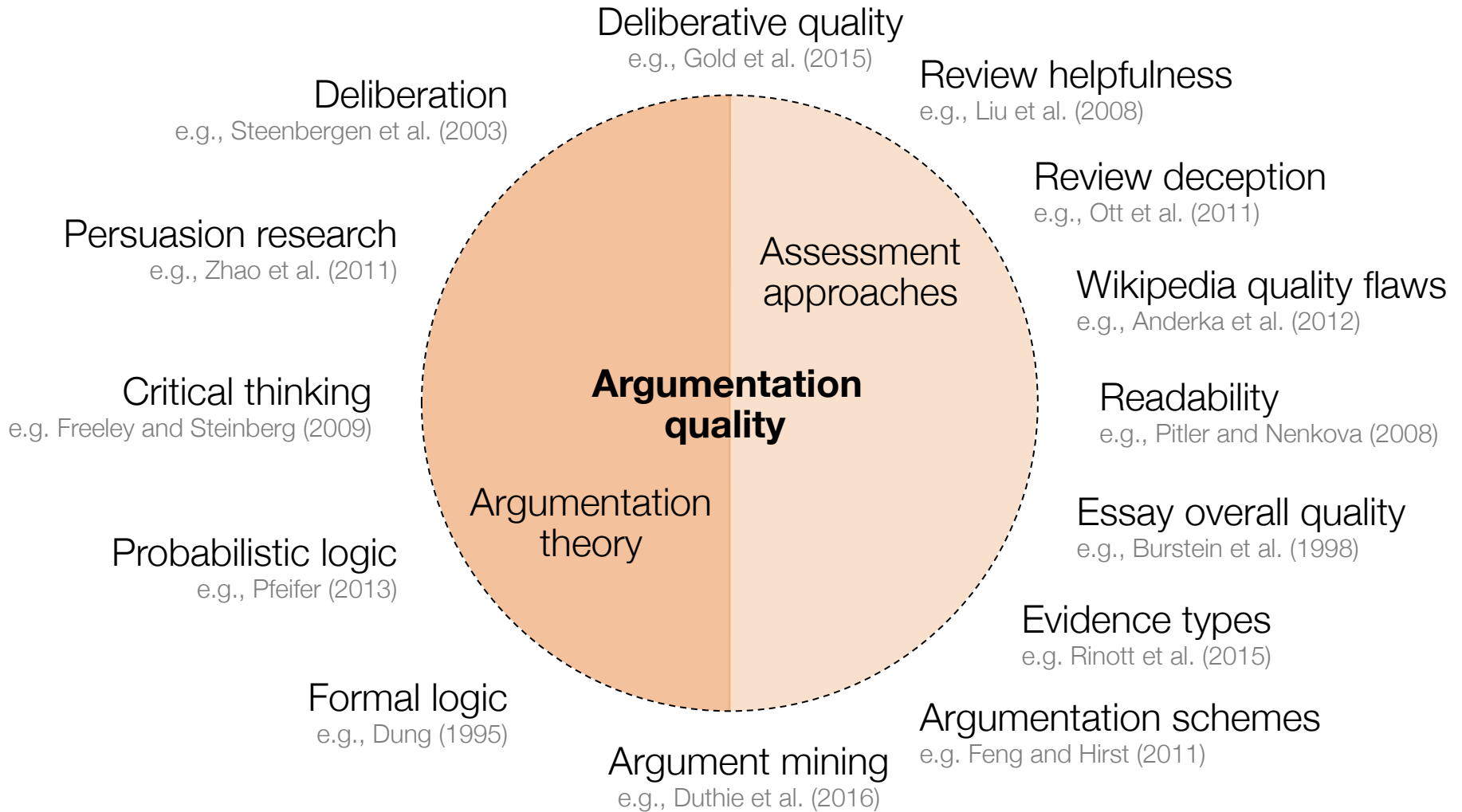
▪ Research questions

- Can we unify the different views of quality?
- Can we provide a common ground for quality assessment?

Core result



Starting point



Survey of existing research

Toulmin (1958) Walton et al. (2008) Cabrio and Villata (2012)
Braunstain et al. (2016) van Eemeren and Grootendorst (2004)
Tindale (2007) Hamblin (1970) Walton (2006) Boltužić and Šnajder (2015)
Rahimi et al. (2014) **Logic** Damer (2009) **Dialectic** Cohen (2011)
Stab and Gurevych (2017) Johnson and Blair (2006) Wachsmuth et al. (2017)
Govier (2010) Blair (2012) **Argumentation quality** Mercier and Sperber (2011)
Freeman (2011) Persing and Ng (2015) van Eemeren (2015) Rahimi et al. (2015)
Persing and Ng (2013) Perelman and Olbrecht-Tyteca (1969) Persing et al. (2010)
Feng et al. (2014) Hoeken (2001) **Rhetoric** Tan et al. (2016) Wei et al. (2016)
Persing and Ng (2014) O’Keefe and Jackson (1995) Zhang et al. (2016)
Park et al. (2015) Aristotle (2007) Habernal and Gurevych (2016)

Three main quality aspects

$$\frac{A \quad A \rightarrow B}{B}$$

Logic

"A dialectical discussion derives its reasonableness from a dual criterion: problem validity and intersubjective validity."

van Eemeren (2015)

$$\frac{A \quad A \rightarrow B}{B}$$

$$\frac{B \rightarrow C}{C}$$

Dialectic



"An argument is cogent if its premises are relevant to its conclusion, individually acceptable, and together sufficient to draw the conclusion."

Blair (2012)

Argumentation
quality

Rhetoric

$$\frac{A \quad A \rightarrow B}{B}$$



"In making a speech, one must study three points: the means of producing persuasion, the style or language to be used, and the proper arrangement of the various parts."

Aristotle (2007)

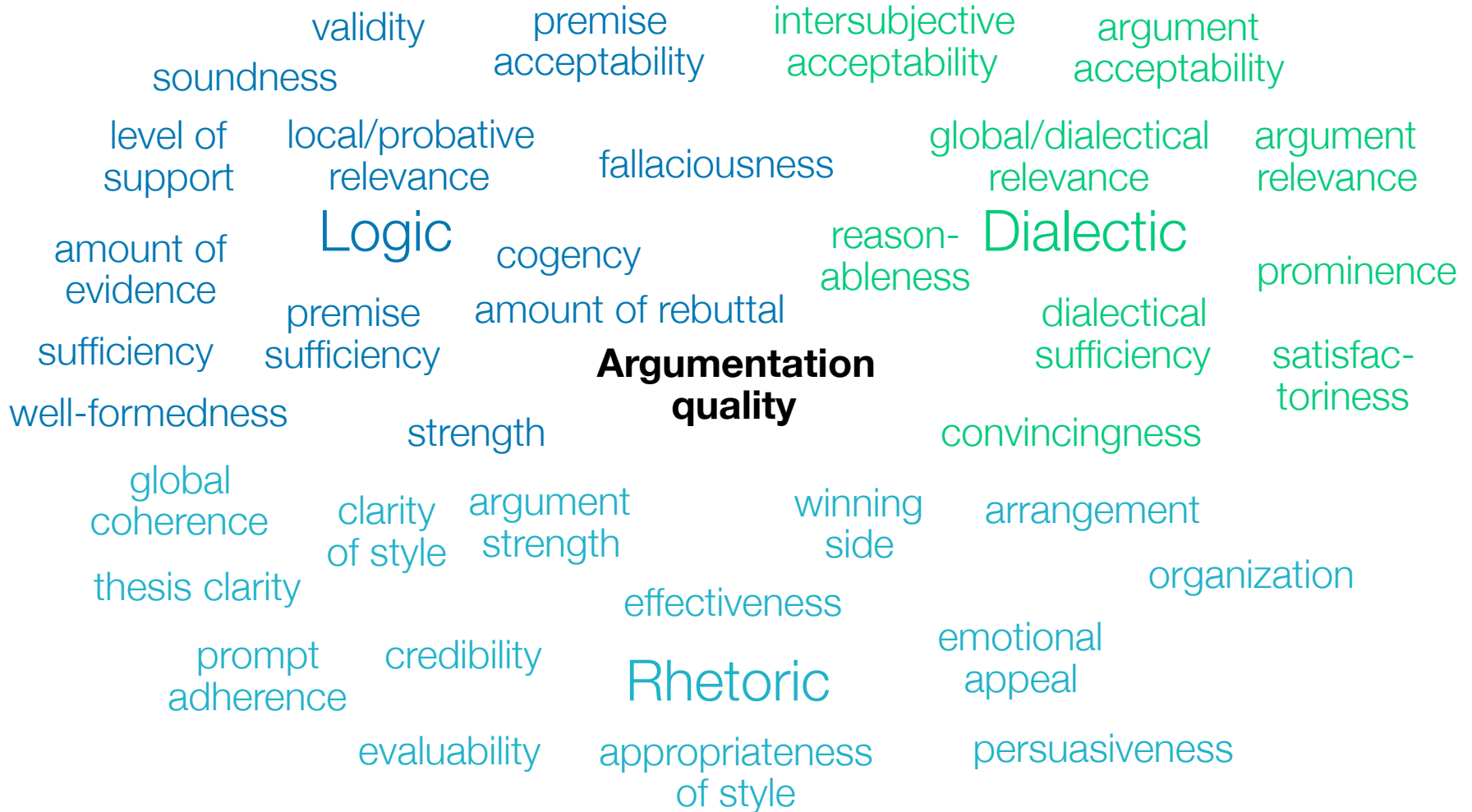
Unification of views

focus on theory

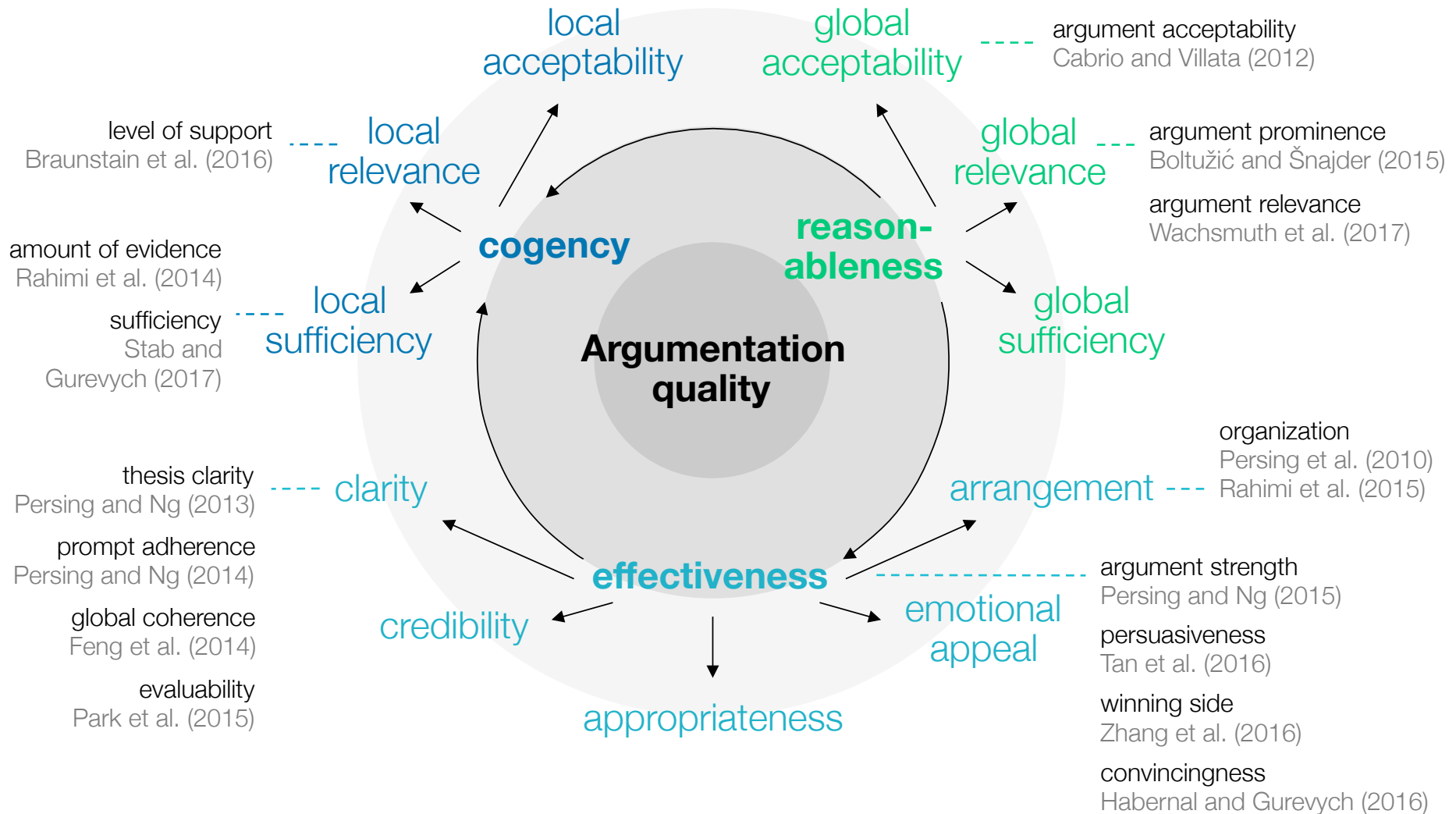
focus on accepted

prefer general

unify names



A taxonomy of computational argumentation quality



The Dagstuhl-15512 ArgQuality corpus

available at
www.arguana.com

- **Corpus based on the taxonomy**
 - 320 debate portal arguments
(Habernal and Gurevych, 2016)
 - 10 per issue/stance pair
 - 3 annotators per argument
 - Score from [1, 3] for all 15 dimensions
- **Agreement**
 - Krippendorff's alpha limited
 - Majority agreement very high
- **Correlations**
 - Overall quality correlates most with reasonableness (.86), cogency (.84), and effectiveness (.81)
 - Several other intuitive correlations

Dimension	Mean	Alpha	Maj.
cogency	1.6	.44	92%
local acceptability	1.9	.46	91%
local relevance	2.3	.47	92%
local sufficiency	1.5	.44	93%
effectiveness	1.4	.45	94%
credibility	1.7	.37	96%
emotional appeal	1.9	.26	94%
clarity	2.1	.35	90%
appropriateness	2.1	.36	88%
arrangement	1.8	.39	93%
reasonableness	1.6	.50	96%
global acceptability	1.9	.44	95%
global relevance	2.0	.42	90%
global sufficiency	1.2	.27	98%
overall quality	1.6	.51	94%

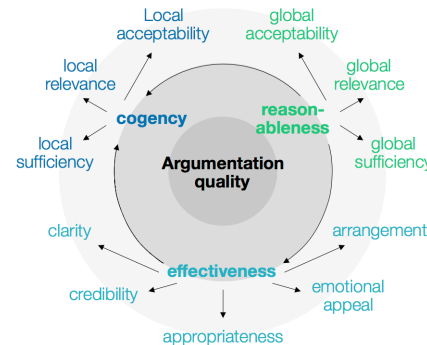
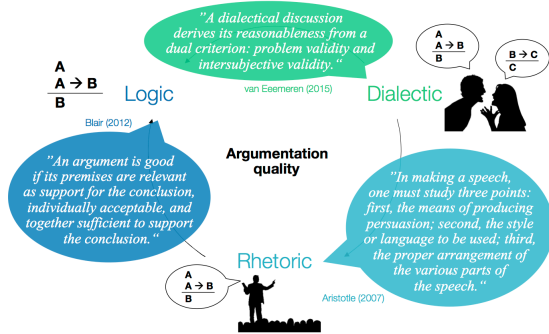
Contributions and outlook

Contributions

Comprehensive survey

Unifying taxonomy

Annotated corpus



Dimension	ϕ	α	Maj.
cogency	1.6	.44	92%
local acceptability	1.9	.46	91%
local relevance	2.3	.47	92%
local sufficiency	1.5	.44	93%
effectiveness	1.4	.45	94%
credibility	1.7	.37	96%
emotional appeal	1.9	.26	94%
clarity	2.1	.35	90%
appropriateness	2.1	.36	88%
arrangement	1.8	.39	93%
reasonableness	1.6	.50	96%
global acceptability	1.9	.44	95%
global relevance	2.0	.42	90%
global sufficiency	1.2	.27	98%
overall quality	1.6	.51	94%

Outlook

Reliable assessment

Target audience

Granularity levels

Theory vs. practice

