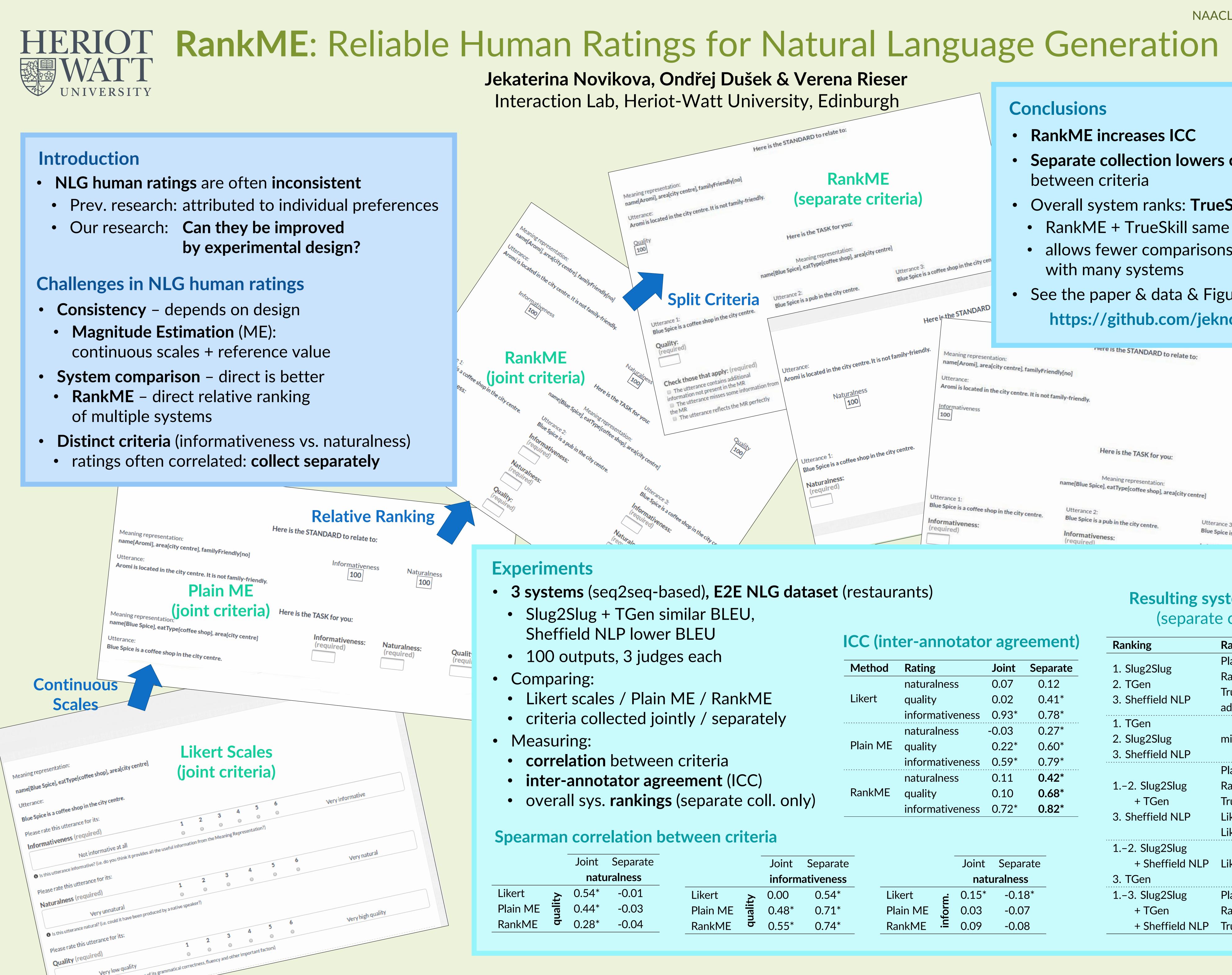


- of multiple systems



			ICC (int	er-anr	otato	Ranking	Rating criterion & method		
ME ately C) coll. only)		Method Likert	naturalness <ert quality<="" th=""><th>Joint 0.07 0.02</th><th><b>Separate</b> 0.12 0.41*</th><th><ol> <li>Slug2Slug</li> <li>TGen</li> <li>Sheffield NLP</li> </ol></th><th colspan="2" rowspan="3"><ul> <li>Plain ME informativeness</li> <li>RankME quality</li> <li>TrueSkill quality</li> <li>added information</li> <li>missing information</li> <li>Plain ME quality</li> <li>RankME informativeness</li> <li>TrueSkill informativeness</li> <li>Likert quality</li> <li>Likert informativeness</li> </ul></th></ert>		Joint 0.07 0.02	<b>Separate</b> 0.12 0.41*	<ol> <li>Slug2Slug</li> <li>TGen</li> <li>Sheffield NLP</li> </ol>	<ul> <li>Plain ME informativeness</li> <li>RankME quality</li> <li>TrueSkill quality</li> <li>added information</li> <li>missing information</li> <li>Plain ME quality</li> <li>RankME informativeness</li> <li>TrueSkill informativeness</li> <li>Likert quality</li> <li>Likert informativeness</li> </ul>	
		informative naturalness Plain ME quality informative		ess	0.93* -0.03 0.22* 0.59*	0.27* 0.60* 0.79* 0.42* 0.68*	1. TGen 2. Slug2Slug 3. Sheffield NLP		
		RankME	informativeness naturalness quality informativeness		0.59* 0.11 0.10 0.72*		1.–2. Slug2Slug + TGen 3. Sheffield NLP		
e	Joint <b>inform</b>	Separate ativeness			Joint na	Sepa turalnes		1.–2. Slug2Slug + Sheffield NLP 3. TGen	Likert naturalness
duaiity	0.00 0.48* 0.55*	0.54* 0.71* 0.74*	Pl	kert ain ME ankME	0.15* 0.03 0.09	-0.1 -0.0 -0.0	)7	1.–3. Slug2Slug + TGen + Sheffield NLP	Plain ME naturalness RankME naturalness TrueSkill naturalness

	Joint	Separate				Joint	Separate
	inform		naturalness				
Ž	0.00	0.54*	Like	ert	'n.	0.15*	-0.18*
quality	0.48*	0.71*		forr	0.03	-0.07	
h	0.55*	0.74*	Rar	hkME	inf	0.09	-0.08

NAACL, New Orleans, June 2018

# • Separate collection lowers correlation

### • Overall system ranks: **TrueSkill efficient**

• RankME + TrueSkill same ranks as RankME • allows fewer comparisons

# • See the paper & data & FigureEight designs: https://github.com/jeknov/RankME/

Here is the TASK for you

Meaning representation name[Blue Spice], eatType[coffee shop], area[city centre]

Utterance 3:

Blue Spice is a coffee shop in the city centre.

## **Resulting system rankings** (separate collection)