



NÚCLEO DE APOIO À PESQUISA EM SOFTWARE LIVRE

Workshop do Núcleo de Apoio à Pesquisa em Software Livre

São Carlos, 16 a 17 de Outubro de 2014.

Mineração de dependências de mudança em repositórios de Software Livre

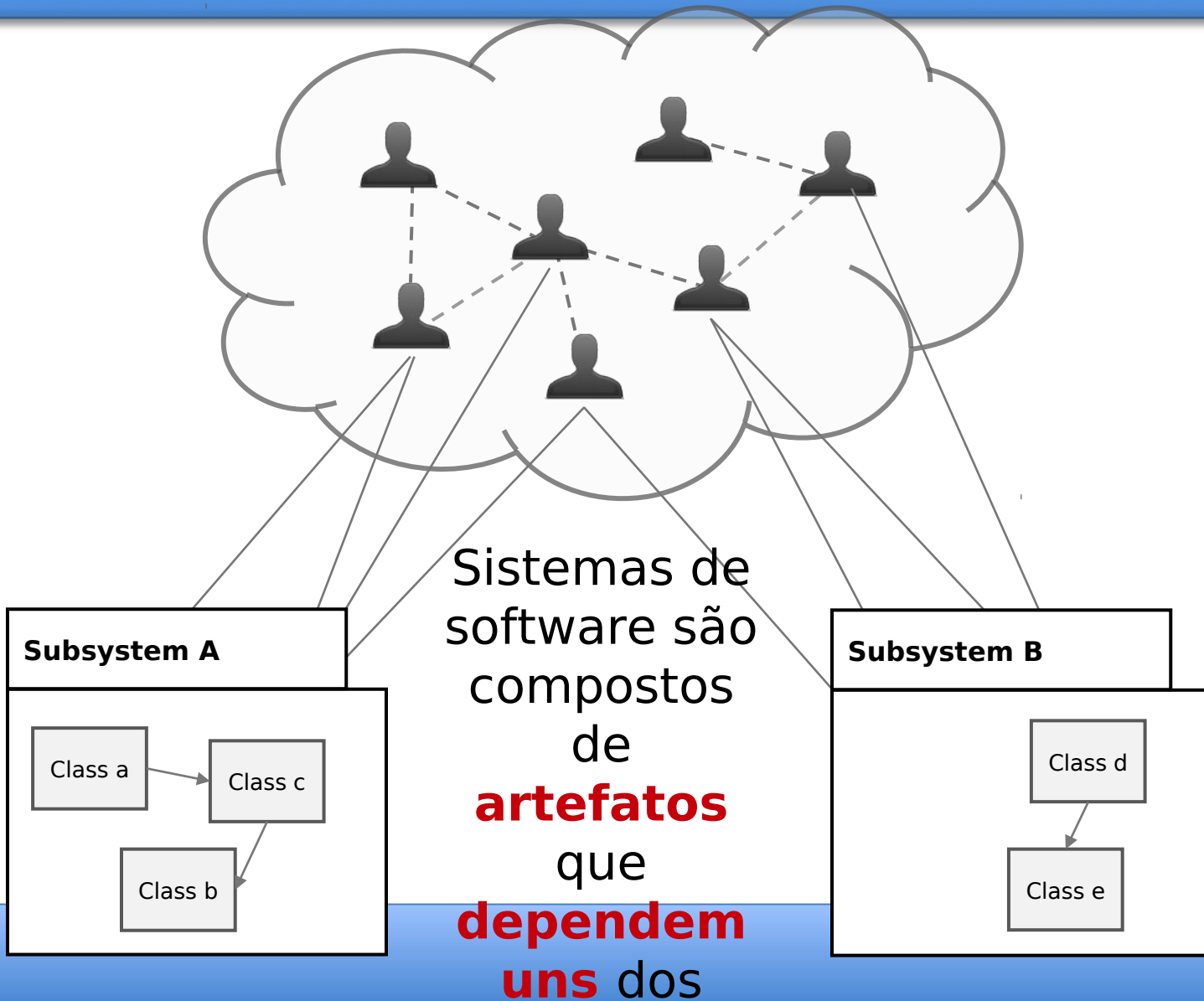
Igor Scaliante Wiese

Universidade Tecnológica Federal do
Paraná / IME-USP

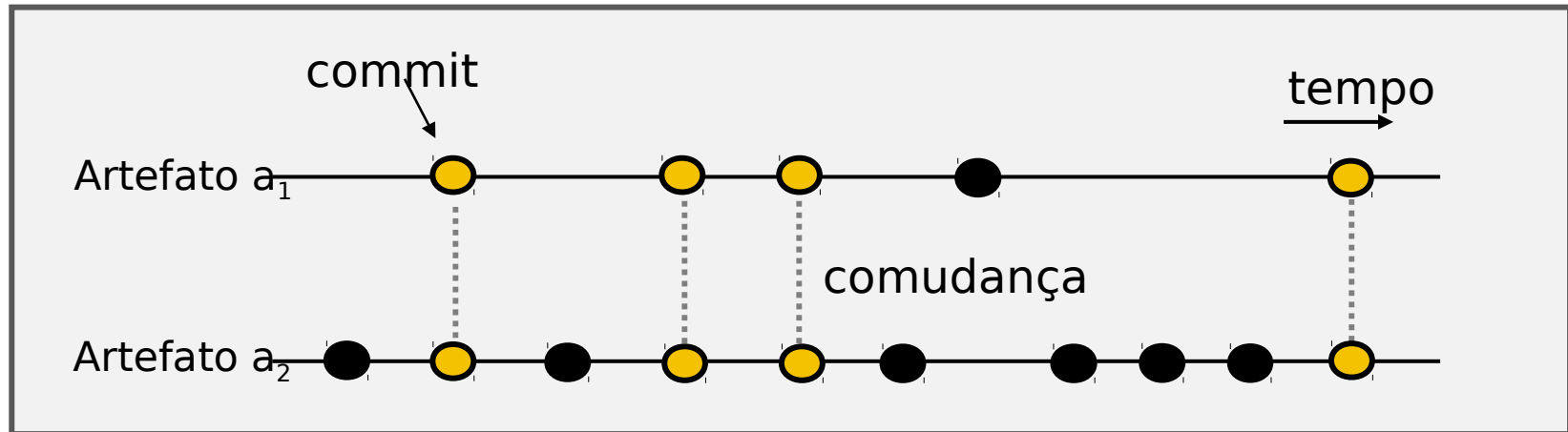
igor.wiese@gmail.com /
igor@utfpr.edu.br



Motivação

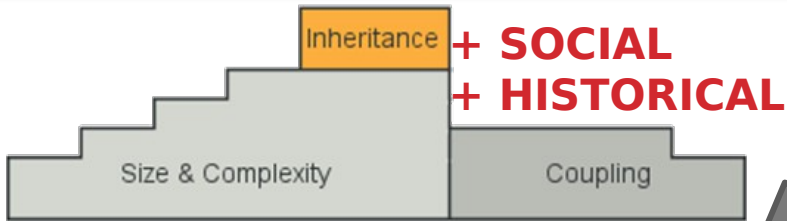


Dependências de Mudança

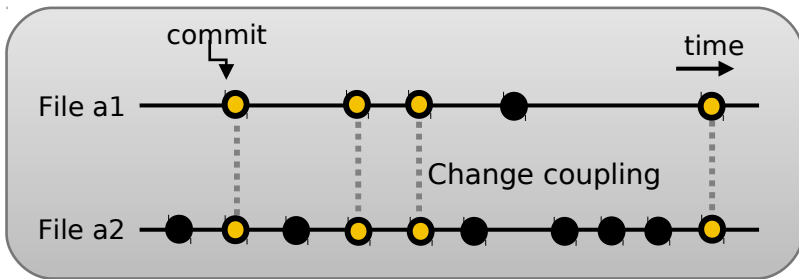


Dependências de mudança são inferidas a partir da análise da evolução histórica. Quanto **mais mudanças conjuntas** dois artefatos têm, **mais acoplados** eles estão.

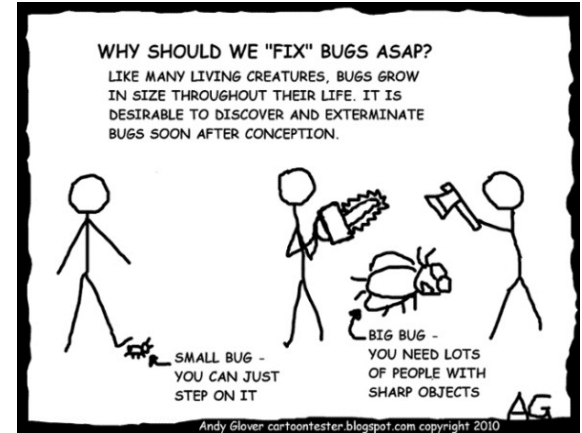
Trabalhos Relacionados



Gustavo Oliva,
Markus Geipel

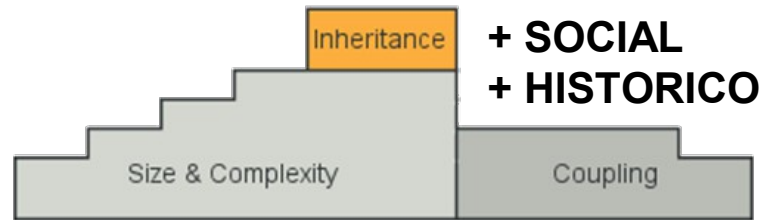


D´ambros - benchmark
Tracy hall - SLR, etc

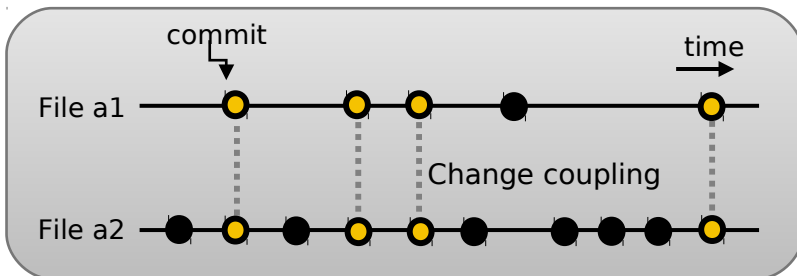


D´ambros - OSS
Kirbas/Ayse Bener - Industrial

Questão de pesquisa central

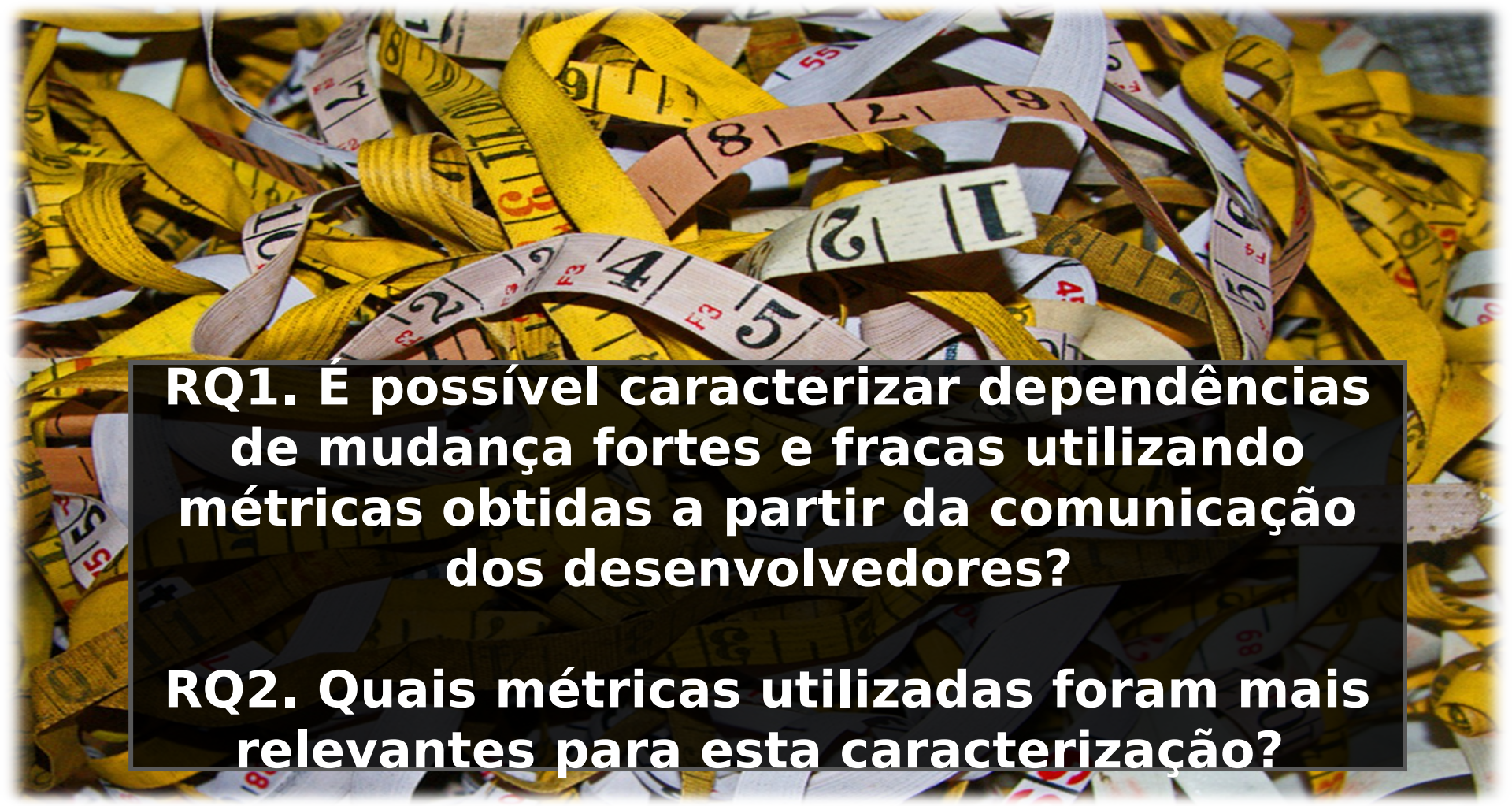


Gustavo Oliva,
Markus Geipel



O quanto as dependências de mudança podem ser explicadas por métricas técnicas, históricas e sociais?

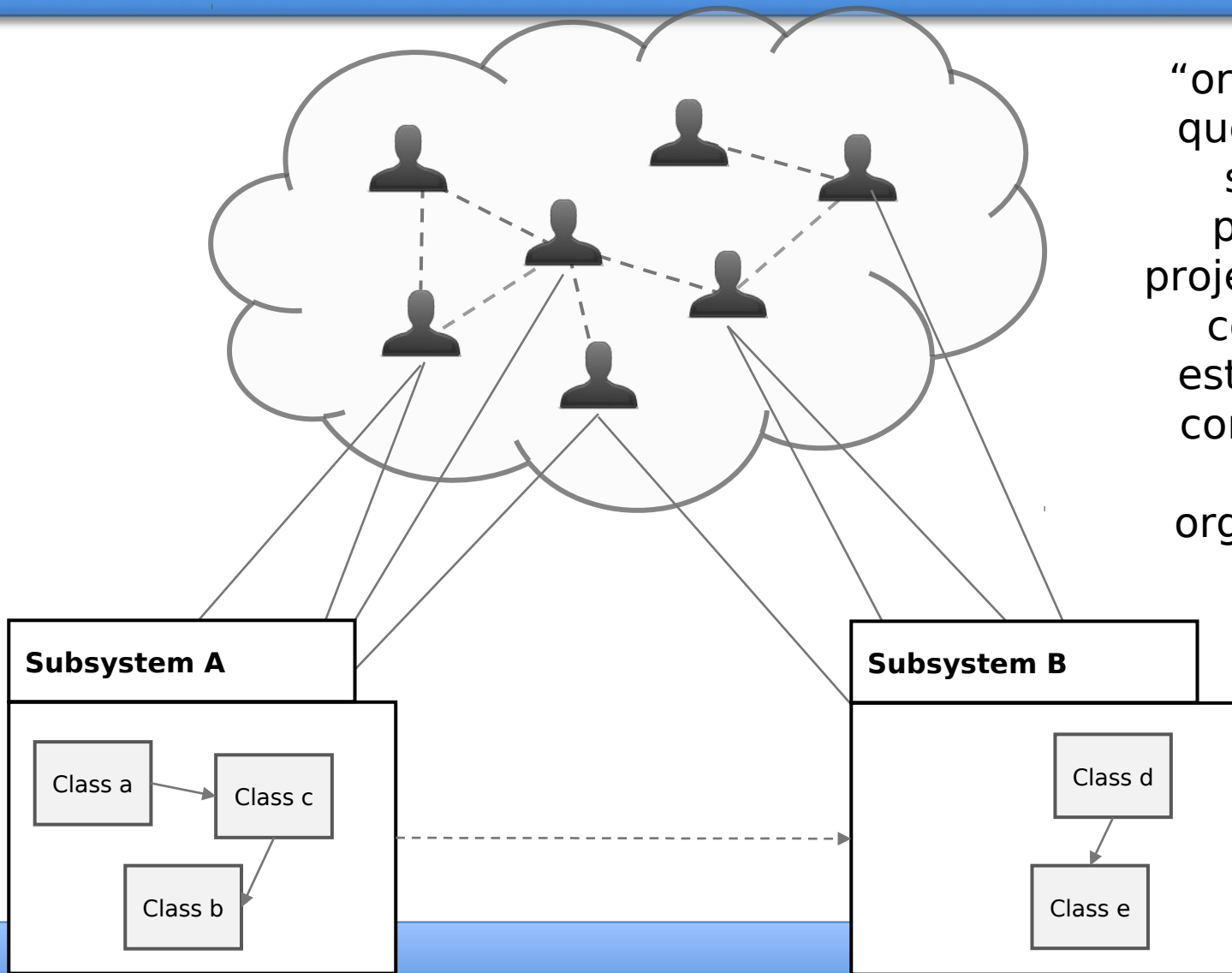
Resultados Preliminares



RQ1. É possível caracterizar dependências de mudança fortes e fracas utilizando métricas obtidas a partir da comunicação dos desenvolvedores?

RQ2. Quais métricas utilizadas foram mais relevantes para esta caracterização?

Conway's Law



“organizações que projetam sistemas produzem projetos que são cópias das estruturas de comunicação desta organização ”



Coleta de Dados

Tabela 1. Sumarização das quatro releases estudadas.

Release	Máximo Suporte	Mediana Suporte	Limiar de Suporte (3º Quartil)	Mediana Confiança	# inst. Fortes	# inst. Fracos	# Pull requests	Total de Instâncias
3.1	14	6	≥ 4	60,00%	873 (45,17%)	1060 (54,83%)	303	1933
3.2	7	2	≥ 2	100,00%	216 (46,87%)	215 (53,13%)	229	431
4.0	14	8	≥ 3	70,17%	513 (39,70%)	778 (60,30%)	1452	1291
4.1	7	2.27	≥ 2	66,66%	339 (16,12%)	1764 (83,88%)	769	339

- 18 métricas sociais.
 - centralidade da rede (*degree*, *betweenness*, *closeness*, e *eigenvector*),
 - medidas de *ego network* (*ego Betweenness*, *ego size*, *ego ties* e *ego density*),
 - medidas de *structural hole* (*efficiency*, *effective size*, *constraint* e *hierarchy*)
 - medidas globais (*size*, *ties*, *density*, *diameter*, número de mensagens e número distintos de desenvolvedores que comentaram).
- 3 métricas do histórico de modificações.
 - *Code Churn* de cada arquivo que forma uma dependência de

É POSSÍVEL CARACTERIZAR DEPENDÊNCIAS DE MUDANÇA FORTES E FRACAS UTILIZANDO MÉTRICAS OBTIDAS A PARTIR DA COMUNICAÇÃO DOS DESENVOLVEDORES?

Tabela 2. Resultados da classificação de dependências de mudança forte e fracas

Release	Algoritmos	% Instâncias corretamente Classificadas	Precisão (forte)	Precisão (fraca)	Sensibilidade (forte)	Sensibilidade (fraca)	AUC
3.1	J48	0,977	0,983	0,974	0,968	0,986	0,977
	Bagging	0,981	0,975	0,973	0,967	0,979	0,995
	KNN	0,980	0,979	0,977	0,971	0,983	0,995
	Naïve Bayes	0,904	0,921	0,892	0,863	0,905	0,943
	JRip	0,967	0,967	0,968	0,961	0,973	0,964
3.2	J48	0,951	0,953	0,949	0,949	0,953	0,951
	Bagging	0,948	0,929	0,971	0,972	0,926	0,986
	KNN	0,974	0,977	0,972	0,972	0,977	0,986
	Naïve Bayes	0,798	0,911	0,734	0,662	0,822	0,875
	JRip	0,951	0,930	0,975	0,977	0,926	0,963
4.0	J48	0,869	0,857	0,876	0,805	0,911	0,875
	Bagging	0,888	0,898	0,883	0,811	0,940	0,947
	KNN	0,874	0,856	0,886	0,823	0,909	0,904
	Naïve Bayes	0,721	0,719	0,723	0,493	0,873	0,713
	JRip	0,793	0,833	0,778	0,602	0,920	0,789
4.1	J48	0,981	0,969	0,984	0,917	0,994	0,976
	Bagging	98,24	0,978	0,983	0,912	0,996	0,986
	KNN	97,86	0,945	0,985	0,92	0,99	0,984
	Naïve Bayes	97,77	0,876	0,911	0,499	0,986	0,805
	JRip	98,11	0,943	0,975	0,920	0,935	0,970

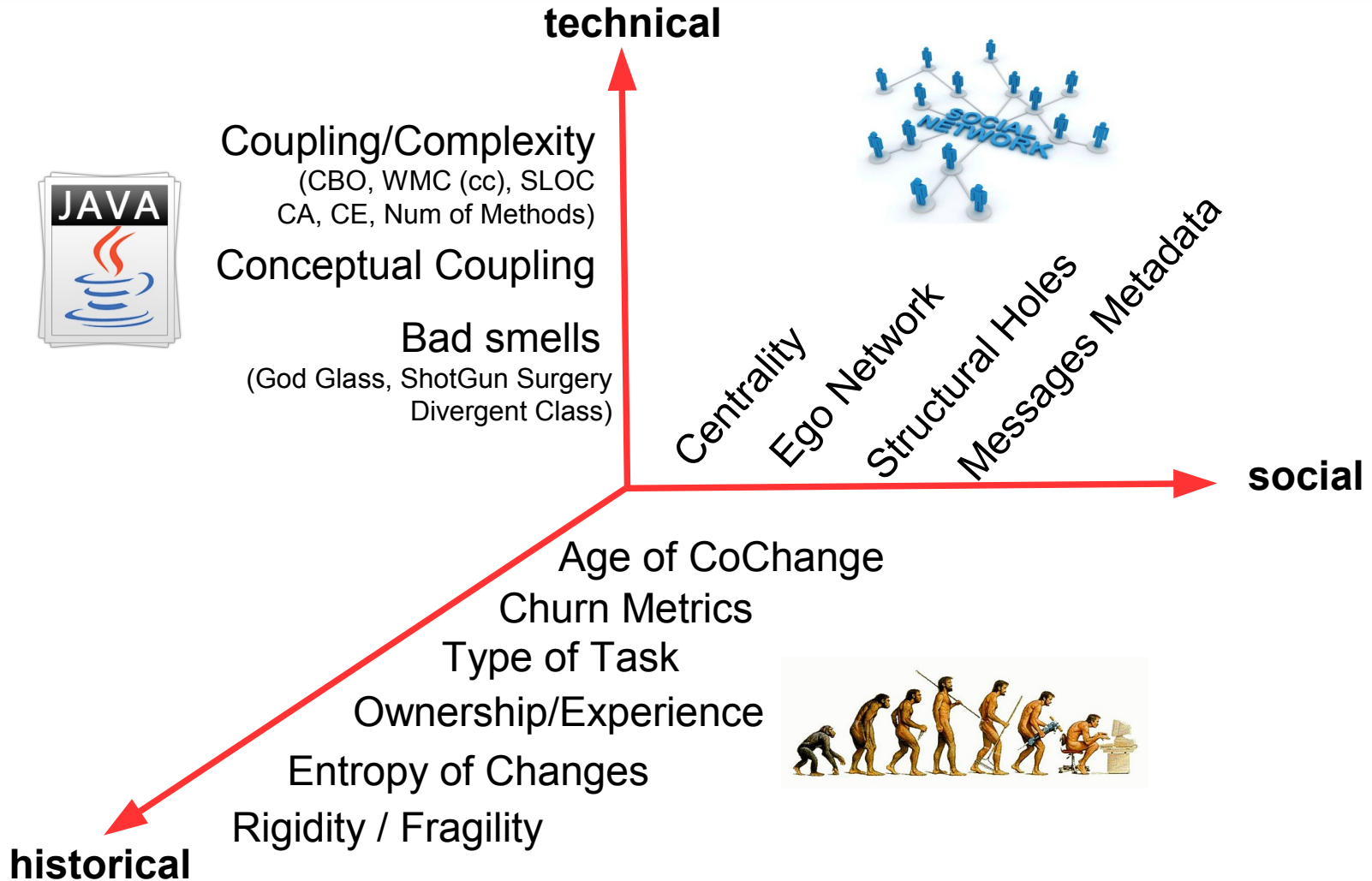
Modelos válidos com AUC > 0.8 em todas as releases

QUAIS MÉTRICAS UTILIZADAS FORAM MAIS RELEVANTES PARA ESTA CARACTERIZAÇÃO?

- Na média, 8,3 métricas (21 no total) foram selecionados pelos 5 algoritmos em 4 releases.
 - A release 3.1 teve média de 9,4 métricas;
 - A release 3.2 teve média de 9,6 métricas;
 - A release 3.2 teve média de 9 métricas
 - A release 3.2 teve média de 5,2 métricas
- Métricas mais selecionadas
 - Densidade e numero de comentários (17 em 20 seleções)
 - *ego Betweenness* e número de desenvolvedores (11 seleções)
 - Restrição, hierarquia e diâmetro (10 seleções)

Evidências que indicam que *o papel de um nó (desenvolvedor) é menos importante do que a estrutura (organização) da rede.*

Próximos passos



OBRIGADO!!!
igor@utfpr.edu.br



<http://lapessc.ime.usp.br/>

