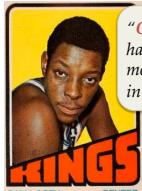


# iCheck: Computationally Combatting “Lies, D—ned Lies, and Statistics”

You Wu<sup>†</sup> Brett Walenz<sup>†</sup> Peggy Li<sup>†</sup> Andrew Shim<sup>†</sup> Seokhyun Song<sup>†</sup> Emre Sonmez<sup>†</sup>

Pankaj K. Agarwal<sup>†</sup> Chengkai Li<sup>‡</sup> Jun Yang<sup>†</sup> Cong Yu<sup>§</sup>

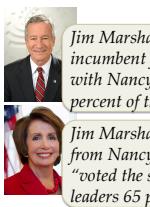
<sup>†</sup>Duke University <sup>‡</sup>Univ. of Texas at Arlington <sup>§</sup>Google Inc.



"Only 10 players in NBA history had more points, more rebounds, and more assists per game than Sam Lacey in their career!"



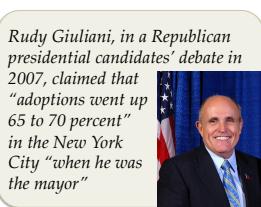
"Jun had 3 SIGMOD papers and 3 ICDE papers in 2006; only seven other researchers have ever beaten this record"



Jim Marshall, a Democratic incumbent from Georgia, voted with Nancy Pelosi "almost 90 percent of the time"



Jim Marshall "is a long way from Nancy Pelosi," as he "voted the same as Republican leaders 65 percent of the time"



Rudy Giuliani, in a Republican presidential candidates' debate in 2007, claimed that "adoptions went up 65 to 70 percent" in the New York City "when he was the mayor"



To check a claim, *tweak* how it manipulates data and see how conclusions differ

## Claim = a parameterized query $q$

- Has a specific setting  $p_0$  of parameters
- Returns a specific answer  $q(p_0)$

*Perturb a claim:* try different  $p$  from the parameter space  $\mathcal{P}$  and compare  $q(p)$  with  $q(p_0)$

- Let  $SR(r; r_0)$  measure the strength of result  $r$  relative to  $r_0$  (+ strengthens it; - weakens it)
- Let  $SP(p; p_0)$  measure the sensibility of parameter setting  $p$  given  $p_0$  (pdf/pmf over  $\mathcal{P}$ —higher means more relevant to checking)

*Query response surface (QRS):* surface over  $\mathcal{P}$  with height at  $p \in \mathcal{P}$  given by  $SR(q(p); q(p_0))$

## Giuliani's claim, parameterized:

- Window aggregate comparison query; result: 66.5
- Three parameters:  $w = 6$  (window length),  $t = 2001$  (end of the second period),  $d = 6$  (distance between two periods)

## How do we check these claims?

**Opportunity:** more and more structured data are available for checking

**Crisis:** Fact-checking is absurdly hard; traditional news media, in decline, can't keep up with investigative reporting needs

**Leverage computing** to "reverse-engineer" vague claims, assess quality beyond correctness, and come up with counterarguments



## Challenge 1: vagueness (by design?)

Adoptions went up 65 to 70 percent in NYC when Giuliani was the mayor

Huh?

Total number of adoptions during 1996-2001 increased by 66.5% compared with 1990-1995

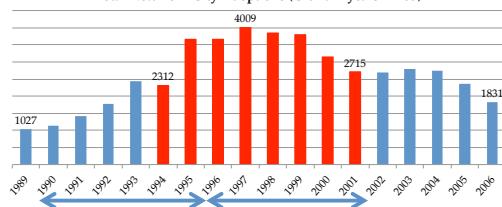
Hmm...

Giuliani was in office 1994-2001

## Challenge 2: beyond correctness

Correct, but still misleading...

Total New York City Adoptions (Giuliani years in red)



## Challenge 3: gimme the punchline

*Counterargument:* Comparing the first and last years of Giuliani's tenure, adoptions increased by only 17%

## Computational challenge: efficient perturbation analysis of database queries

- Exploit knowledge of queries
- Exploit properties of sensibility functions

[tinyurl.com/icheckuclaim](http://tinyurl.com/icheckuclaim)

This is just the tip of an iceberg: join us in exploring endless possibilities for applying computing to fact-checking, journalism, and beyond!

