



# Supervised vs. Unsupervised Automated Categorization

Joseph Busch, Principal Analyst

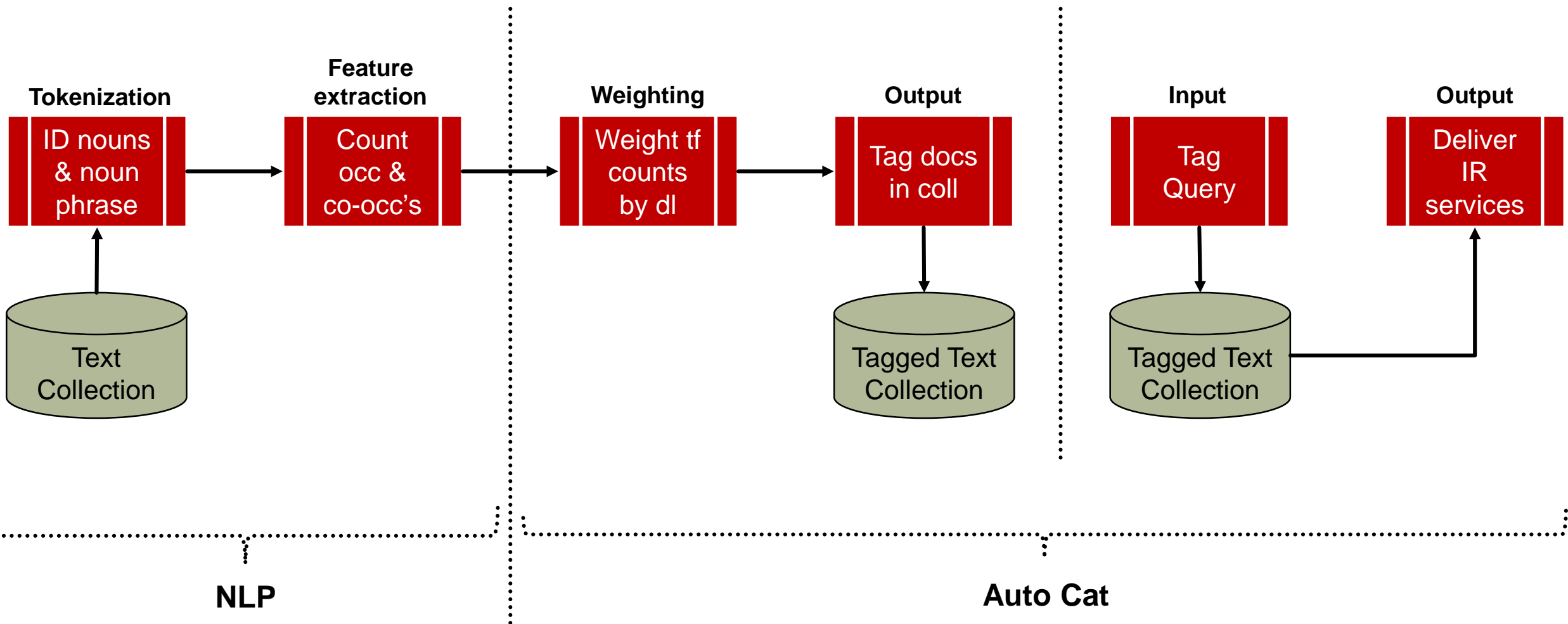
# Agenda

- ❖ How to tag content using text analytics
- ❖ How to choose appropriate target collections
- ❖ How to evaluate text analytics accuracy
- ❖ RWJF case study

# Categorization goals

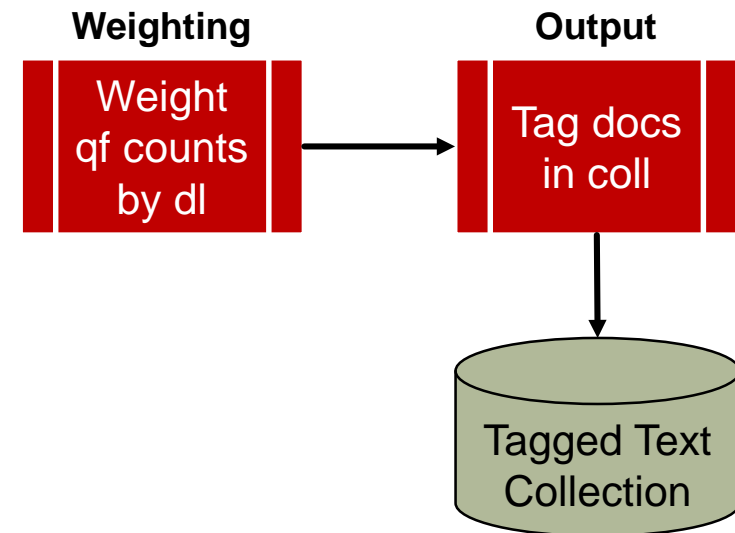
- ❖ Make sense
- ❖ Provide clear perception
- ❖ Is trustworthy and as unbiased as possible

# Natural language processing enables unsupervised automated categorization



# Individual terms can be replaced by pre-defined queries chosen and constructed by editors

- ❖ tf becomes qf
- ❖ Simple pre-defined query for “Affordable Care Act”
  - ("Affordable Care Act" OR "ACA" OR Obamacare)
- ❖ Complex pre-defined query for “Medicare Navigator”
  - ((navigator\* OR assistor\* OR assister\*) NEAR ("health insurance" OR Medicare OR Medicaid) NEAR enroll\*))



# Agenda

- ❖ How to tag content using text analytics
- ❖ How to choose appropriate target collections
- ❖ How to evaluate text analytics accuracy
- ❖ RWJF case study

# How reliable does a categorizer need to be?

It depends on the categorization target

**Documents**



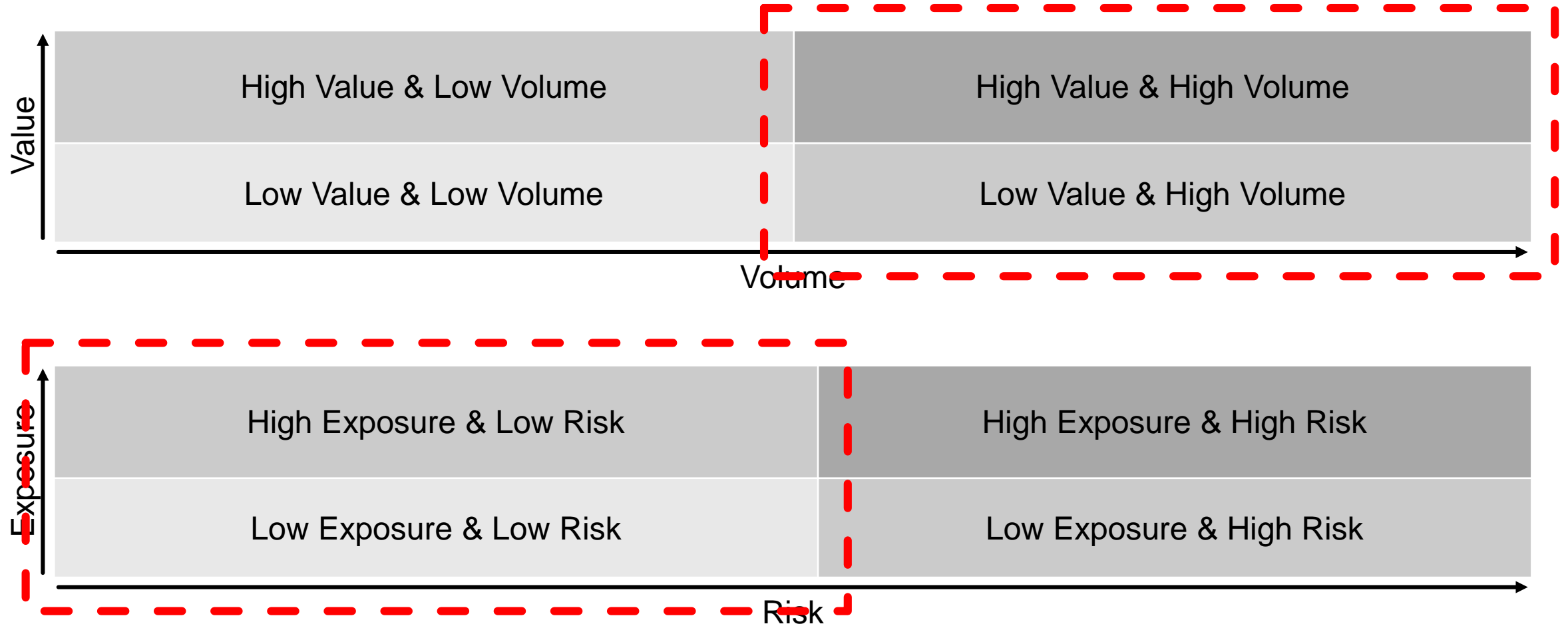
**VS.**

**Projects**



# Productivity vs risk models

## Low vs high error tolerance





# Agenda

- ❖ How to tag content using text analytics
- ❖ How to choose appropriate target collections
- ❖ How to evaluate text analytics accuracy
- ❖ RWJF case study

# How to evaluate text analytics accuracy

- ❖ Build a prototype to test feasibility of auto-categorization.
  - Test selected categories individually
    - Is the recall and precision optimized on a category basis?
  - Test selected categories together
    - Can the categorizer appropriately distinguish between categories?
    - Are multiple categories identified when that is appropriate?
    - Are there relevance measures that can be meaningfully used to distinguish among multiple identified categories?
- ❖ Build production-quality categorizers.
  - Test all categories individually to optimize recall and precision
  - Test all categories together to ensure they can be appropriately distinguished

# Agenda

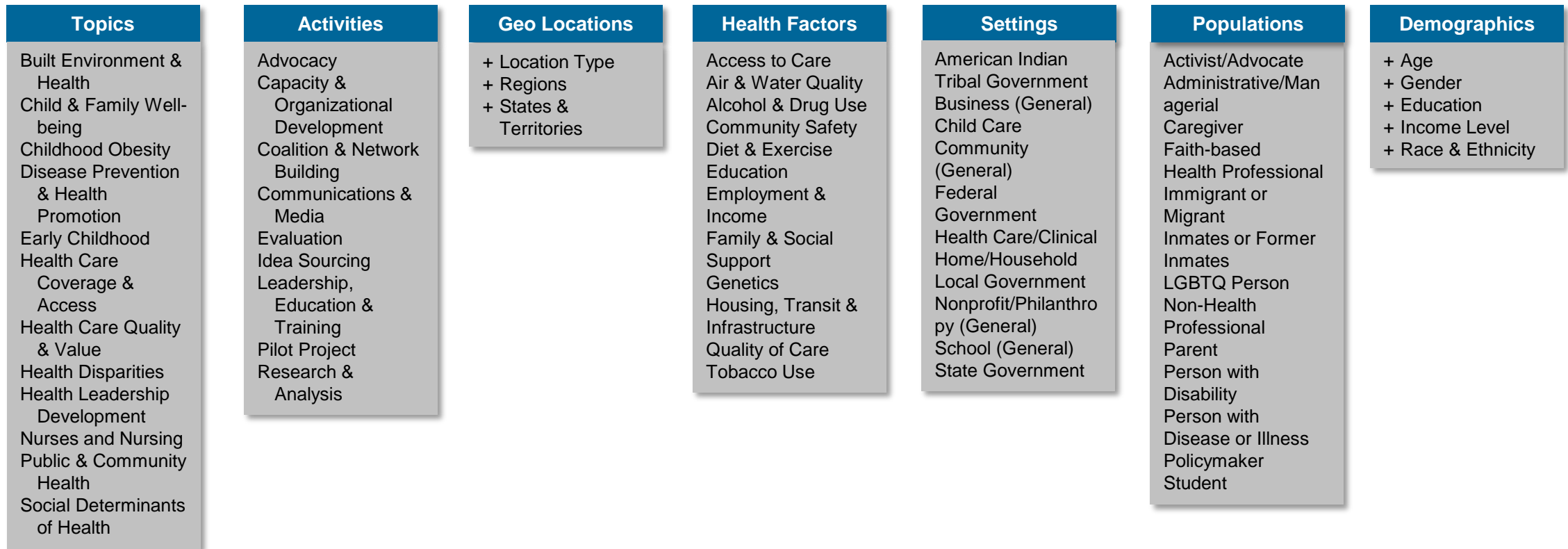
- ❖ How to tag content using text analytics
- ❖ How to choose appropriate target collections
- ❖ How to evaluate text analytics accuracy
- ❖ RWJF case study

# Former “A-Z” RWJF Taxonomy for their Project Information Management System (PIMS)

<ul style="list-style-type: none"> <li>Access and barriers to care</li> <li>Accountable care organizations</li> <li>Accreditation</li> <li>Addiction and substance abuse</li> <li>Administration</li> <li>Advanced practice nurses</li> <li>Affordable Care Act (ACA)</li> <li>Aging</li> <li>Alcoholism</li> <li>Allied health</li> <li>Alternative and spiritual care</li> <li>Ambulatory care</li> <li>American Indian tribal government</li> <li>Antibiotic resistance</li> <li>Asthma</li> <li>At-risk and vulnerable people</li> <li>Barriers to care: cultural, gender and racial</li> <li>Barriers to care: financial</li> <li>Barriers to care: language and literacy</li> <li>Barriers to care: logistics and transportation</li> <li>Behavior change</li> <li>Behavioral economics</li> <li>Behavioral/mental health</li> <li>Benchmarks and best practices</li> <li>Built environment</li> <li>Bundled payments</li> <li>Business</li> <li>Cancer</li> <li>Care transitions</li> <li>Careers</li> <li>Case management</li> <li>Child abuse and neglect</li> <li>Child care centers</li> <li>Child welfare</li> <li>Childhood obesity</li> <li>Children’s health insurance program</li> <li>Chronic disease management</li> <li>Chronic illness</li> <li>Clean air laws</li> <li>Clinical care</li> <li>Clinical research</li> <li>Colleges and universities</li> <li>Community benefit</li> <li>Community development</li> <li>Community outreach</li> </ul>	<ul style="list-style-type: none"> <li>Community violence</li> <li>Community-based care</li> <li>Competitions</li> <li>Competitive foods</li> <li>Complete streets</li> <li>Conferences</li> <li>Consumer engagement</li> <li>Continuing education</li> <li>Coordinated care</li> <li>Cost of care</li> <li>Crowd sourcing</li> <li>Cultural competence</li> <li>Data</li> <li>Dental schools</li> <li>Dentists</li> <li>Developmental disabilities</li> <li>Diabetes</li> <li>Disparities</li> <li>Disruptive innovations</li> <li>Diverse communities</li> <li>Diversity</li> <li>Early childhood</li> <li>Early intervention</li> <li>Education and training</li> <li>Education level</li> <li>E-health</li> <li>Eligibility and enrollment</li> <li>Emergency care</li> <li>Emergency preparedness and response</li> <li>Employers</li> <li>Employer-sponsored insurance</li> <li>End-of-life care</li> <li>Environmental health</li> <li>Epidemiology</li> <li>Episodes of care</li> <li>Ethics/bioethics</li> <li>Evaluation</li> <li>Evidence-based</li> <li>Families</li> <li>Family caregiving</li> <li>Federal government</li> <li>Food access</li> <li>Food safety</li> <li>Foster care</li> <li>Frontline workforce</li> <li>Genetics</li> <li>Health care delivery system</li> </ul>	<ul style="list-style-type: none"> <li>Health care markets</li> <li>Health education</li> <li>Health games</li> <li>Health impact assessments</li> <li>Health insurance</li> <li>Health insurance exchanges</li> <li>Health IT</li> <li>Health plans</li> <li>Health policy</li> <li>Health promotion and disease prevention</li> <li>Health records/electronic health records</li> <li>Health reform</li> <li>Healthy communities</li> <li>Heart disease</li> <li>Home health care</li> <li>Home visiting</li> <li>Homeless</li> <li>Hospital readmissions</li> <li>Hospital-acquired infections</li> <li>Hospitals</li> <li>Housing/public housing</li> <li>Immigrants and refugees</li> <li>Immunizations</li> <li>Independent living or self-determination</li> <li>Individual insurance</li> <li>Industry regulation</li> <li>Infectious diseases</li> <li>Informatics</li> <li>Information technology</li> <li>Injury</li> <li>Inmates and former inmates</li> <li>Inpatient care</li> <li>Institutional design</li> <li>Insurance cost</li> <li>Interpersonal violence</li> <li>Interprofessional collaboration</li> <li>Jails and prisons</li> <li>Job satisfaction</li> <li>Jobs</li> <li>Joint use facilities</li> <li>Journalists</li> <li>Language services</li> <li>Leadership development</li> <li>Legal systems and issues</li> <li>Lesbian/gay/bisexual/transgender</li> <li>Local government</li> <li>Long-term care</li> </ul>	<ul style="list-style-type: none"> <li>Low-birthweight infants</li> <li>Managed care organizations</li> <li>Marketing</li> <li>Medicaid</li> <li>Medical errors</li> <li>Medical homes</li> <li>Medical malpractice</li> <li>Medical practices</li> <li>Medical schools</li> <li>Medical students and residents</li> <li>Medical technology</li> <li>Medical, dental and nursing workforce</li> <li>Medically underserved areas</li> <li>Medicare</li> <li>Mentoring</li> <li>Military/veterans</li> <li>Mobile health/mhealth</li> <li>Models</li> <li>Networks</li> <li>Nurse executives</li> <li>Nurse midwives</li> <li>Nurse practitioners</li> <li>Nurses</li> <li>Nursing homes</li> <li>Nursing schools</li> <li>Nutrition</li> <li>Nutrition policy</li> <li>Obesity policy</li> <li>Open sourcing</li> <li>Oral health</li> <li>Out-of-school time</li> <li>Overweight</li> <li>Palliative care</li> <li>Parks or playgrounds</li> <li>Patient safety and outcomes</li> <li>Patient satisfaction</li> <li>Patient-centered care</li> <li>Patients</li> <li>Pay for performance</li> <li>Payment reform</li> <li>Pediatric care</li> <li>Performance standards and measurement</li> <li>Philanthropy</li> <li>Physical activity</li> <li>Physical activity policy</li> <li>Physical education</li> </ul>	<ul style="list-style-type: none"> <li>Physician assistants</li> <li>Physicians</li> <li>Poor and economically disadvantaged</li> <li>Practice guidelines</li> <li>Pregnancy</li> <li>Prenatal and neonatal care</li> <li>Pre-schools</li> <li>Prescription drugs</li> <li>Prevention</li> <li>Preventive care</li> <li>Primary care</li> <li>Primary care/generalist physicians</li> <li>Prisoner reentry</li> <li>Privacy and confidentiality</li> <li>Provider incentives</li> <li>Public health</li> <li>Public health agencies</li> <li>Public health law</li> <li>Public health professionals</li> <li>Public health schools</li> <li>Public health system and finance</li> <li>Public policy</li> <li>Public-private partnerships</li> <li>Quality of care</li> <li>Recruitment and retention</li> <li>Registered nurses</li> <li>Religious congregations</li> <li>Research</li> <li>Research networks</li> <li>Risky behavior</li> <li>RWJF news</li> <li>Safe routes to school</li> <li>Safety nets</li> <li>Scholars and fellows</li> <li>School foods</li> <li>School snacks</li> <li>School/district policy</li> <li>School-based health centers</li> <li>Schools K-12</li> <li>Scope of practice</li> <li>Screening</li> <li>Self-care</li> <li>Shared decisionmaking</li> <li>Shelters</li> <li>Shortage of medical or nursing personnel</li> <li>Small businesses</li> <li>Social and emotional learning</li> </ul>	<ul style="list-style-type: none"> <li>Social determinants of health</li> <li>Social isolation</li> <li>Social marketing</li> <li>Social sciences</li> <li>Social support services</li> <li>Specialist physicians</li> <li>Specialist physicians</li> <li>State government</li> <li>Substance abuse treatment</li> <li>Substance use</li> <li>Substance use</li> <li>Sugary beverages</li> <li>Sugary beverages</li> <li>Supportive housing</li> <li>Tax policy</li> <li>Tax policy</li> <li>Telemedicine and telehealth</li> <li>Tobacco</li> <li>Tobacco</li> <li>Tobacco cessation</li> <li>Tobacco cessation</li> <li>Tobacco control</li> <li>Transdisciplinary</li> <li>Transparency/public reporting</li> <li>Transportation</li> <li>Transportation</li> <li>Transportation policy</li> <li>Transportation policy</li> <li>Transportation policy</li> <li>Transportation policy</li> <li>Underinsured</li> <li>Underserved populations</li> <li>Underserved populations</li> <li>Underserved populations</li> <li>Uninsured</li> <li>Value-based purchasing</li> <li>Walking and biking</li> <li>Walking and biking</li> <li>Walking and biking</li> <li>Work environment</li> <li>Workflow</li> <li>Workforce issues</li> <li>Workforce supply and demand</li> <li>Workforce supply and demand</li> <li>Youth development</li> </ul>
--	---	--	---	--	---



# New RWJF Taxonomy model for their Project Information Management System (PIMS)



**Problem:** How to convert legacy collection from as-is to the new taxonomy model?

- ❖ Option 1: Mapping from the as-is to the new.
- ❖ Option 2: Automated categorization
- ❖ Option 3: Combination of option 1 and 2

# Option 1: Mapping

<b>CODING _DEF_ID</b>	<b>PIMS Topics</b>		<b>New Facet</b>	<b>New Catgeory</b>
350	Access and barriers to care		Topics	Health Care Coverage & Access
1707	Accountable care organizations		Topics	Health Care Quality & Value
386	Accreditation		Topics	Public & Community Health
508	Addiction and substance abuse		Health Fac	Alcohol & Drug Use
508	Addiction and substance abuse		Topics	Disease Prevention & Health Promotion
670	Administration		Population	Administrative/Managerial
670	Administration		Population	Health Professional
670	Administration		Topics	Health Leadership Development
470	Advanced practice nurses		Population	Health Professional
470	Advanced practice nurses		Topics	Nurses & Nursing
1709	Affordable Care Act (ACA)		Topics	Health Care Coverage & Access
1667	Aging		Demograp	Older Adults
509	Alcoholism		Health Fac	Alcohol & Drug Use
509	Alcoholism		Population	Person with Disease or Illness
509	Alcoholism		Topics	Disease Prevention & Health Promotion
478	Allied health		Population	Health Professional
478	Allied health		Topics	Health Leadership Development
406	Alternative and spiritual care		Settings	Health Care/Clinical
395	Ambulatory care		Settings	Health Care/Clinical
395	Ambulatory care		Topics	Health Care Quality & Value
306	American Indian tribal government		Settings	American Indian Tribal Government

## Option 2-Step 1: Build a prototype to test feasibility of auto-categorization

Can we build some categorizers that accurately detect Topics?

- ❖ Collected 200 categorized examples for each Topic to be tested.
- ❖ Divided test collection in half – 100 to use to build the categorizer, and 100 to use in a blind test.
- ❖ Reviewed 100 relevant items and identify the words and phrases that provide a context for the topic.
- ❖ Noted any named entities (people, organizations, events, laws, etc.) that are closely associated with the topic.
- ❖ Consolidated the terms – identify duplicates, synonyms, as well as any concepts that you want to combine even if they are not synonyms.
- ❖ Wrote a query for each term.
- ❖ Combined the terms into a single nested query.

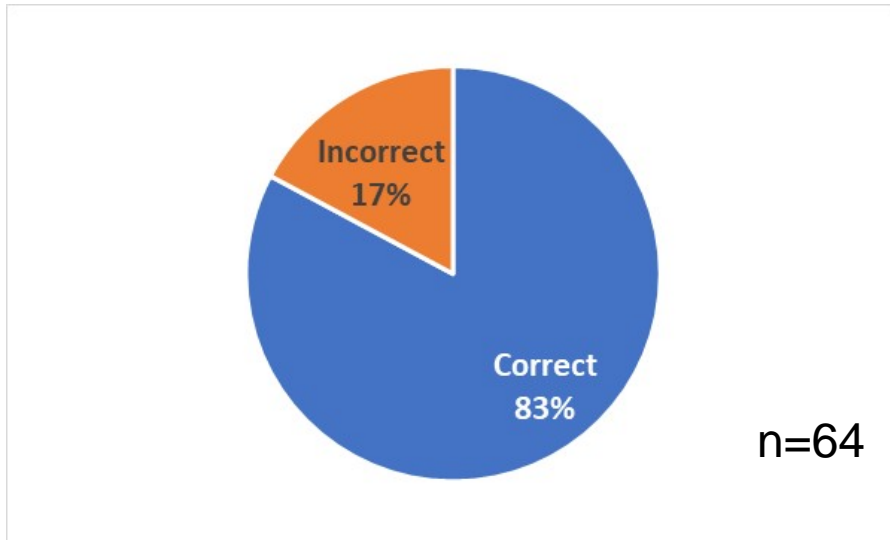
## Pre-defined query for “Health Care Coverage & Access” topic

((("health insurance coverage" OR "health coverage") OR ("healthcare reform" OR "health care reform") OR ("Better Care Reconciliation Act" OR BCRA) OR ("American Health Care Act" OR AHCA) OR ("Affordable Care Act" OR "ACA" OR Obamacare) OR ((Medicare OR Medicaid) NEAR/5 (spend\* OR cover\* OR expan\*)) OR ("health insurance exchange" OR "HIE") OR ("health insurance" NEAR/5 marketplace\*) OR ("federal\* facilitated marketplace\*" NEAR/10 "health insurance") OR ("federal\* run marketplace\*" NEAR/10 "health insurance") OR ((state NEAR/5 marketplace\*) NEAR/10 "health insurance") OR ("small business marketplace\*" NEAR/10 "health insurance") OR ("small-business marketplace\*" NEAR/10 "health insurance") OR (("small business" NEAR/5 exchange\*) NEAR/10 "health insurance") OR (("high-risk" OR "high risk") NEAR/10 "health insurance") OR (uninsured NEAR/5 (veteran\* OR child\* OR adult\* OR people OR kid\* OR citizen\*)) OR (("pre-existing condition\*" OR "preexisting condition\*") NEAR/10 "health insurance") OR "health insurance rate\*" OR ((cost\* OR rate\* OR payment\*) NEAR/10 "health insurance") OR ("health insurance" NEAR/10 "tax credit\*") OR ((healthcare OR "health care") NEAR/5 spending) OR ((healthcare OR "health care") NEAR/5 utilization) OR (("high-deductible" OR "high deductible") NEAR/10 "health insurance") OR (("mental health" OR "substance abuse") NEAR/10 "health insurance") OR ("provider network\*" NEAR/10 "health insurance") OR (("in-network" OR "out-of-network") NEAR/10 "health insurance") OR ((PPO\* OR HMO\*) NEAR/5 (marketplace\* OR plan\* OR provider\*)) OR ("health insurance" NEAR/10 (enroll\* OR "re-enroll\*" OR renew\* OR "open-enrollment" OR "open enrollment")) OR ((navigator\* OR assistor\* OR assister\*) NEAR/10 (("health insurance" OR Medicare OR Medicaid) NEAR/5 enroll\*)) OR ("CHIP" OR "Children’s Health Insurance Program") OR ("individual mandate" NEAR/10 "health insurance") OR "employer-sponsored insurance" OR ((employer OR employee) NEAR/10 "health insurance"))

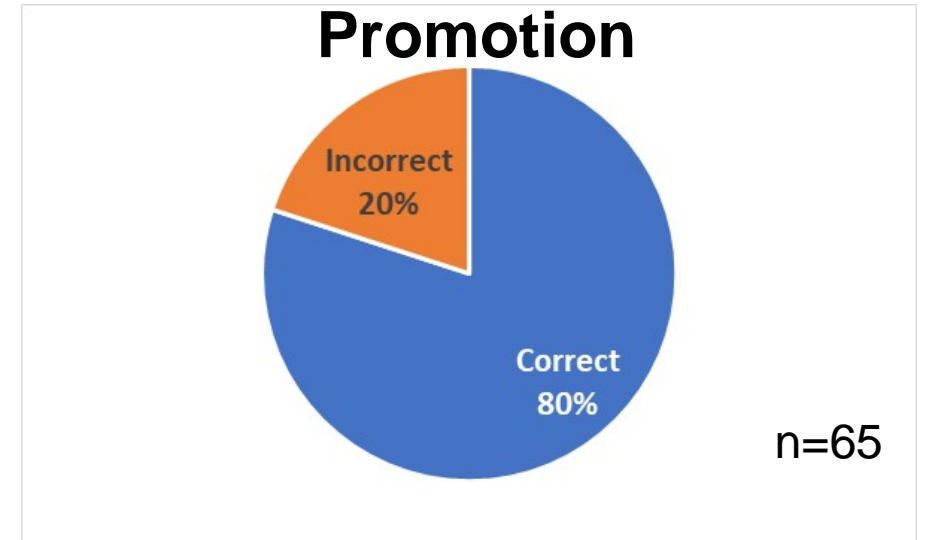


# Trial results for each pre-defined category

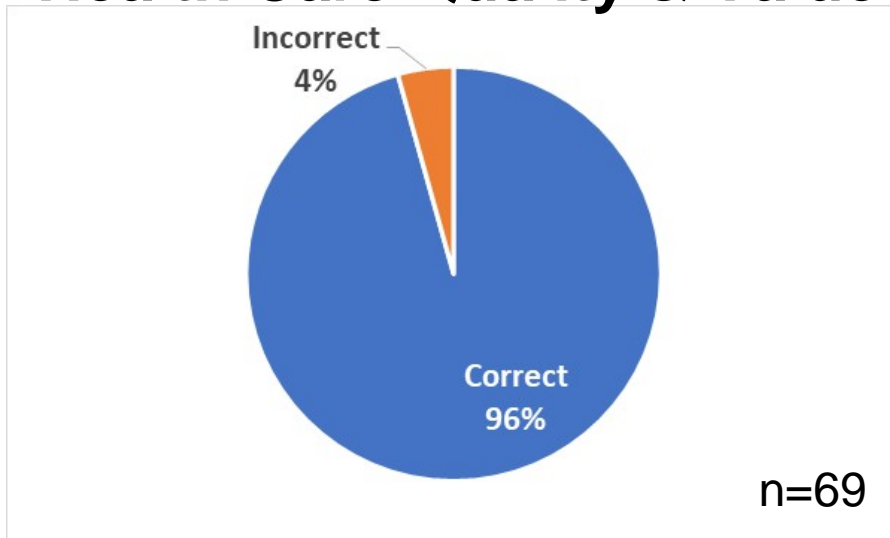
## Childhood Obesity



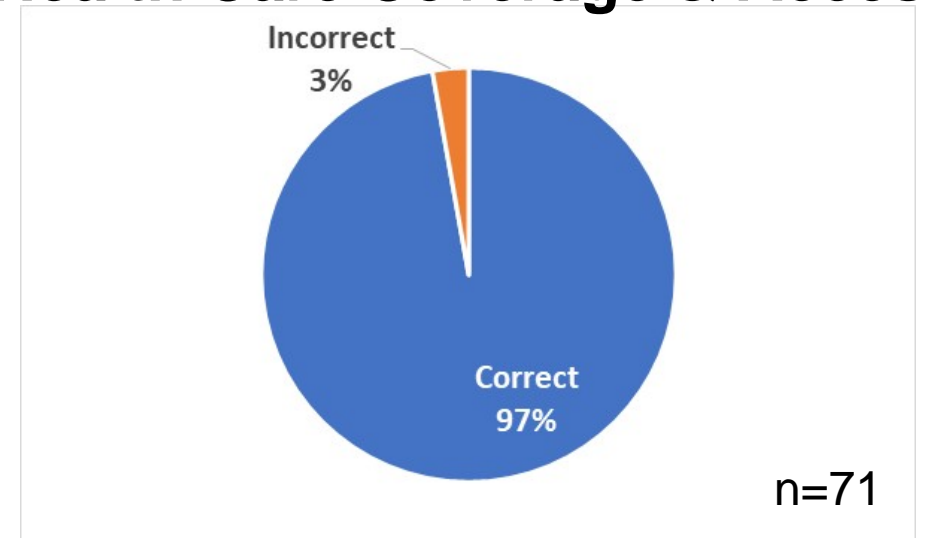
## Disease Prevention & Health Promotion



## Health Care Quality & Value



## Health Care Coverage & Access



## Option 2-Step 2: Build production-quality categorizers for all 12 topics

- ❖ For each of the 12 RWJF Topics:
  - Gather a test collection of grant summaries
  - Build and test a rule – a Boolean query with proximity operators to classify into an RWJF Topic
    - Start with the RWJF Topic string.
    - Add synonyms.
    - Add legacy PIMS Topics that were mapped to the 12 RWJF Topics and synonyms.
    - Review the test collection of grant titles & summaries tagged for that RWJF Topic. (A sample size of 100 should suffice). Identify distinctive words and phrases.
    - Group Boolean search strings together by each discrete topic as much as possible.
    - Review and revise.
  - Test a categorizer for an RWJF Topic. (12 categorizers; 12 tests)
    - Test **recall** – Are **all** the relevant grants being identified by the query, 70-80% of the time? If not, revise query and re-test.
    - Test **precision** – Are **only** the relevant grants being identified by the query, 70-80% of the time? If not, revise query and re-test.
- ❖ Consolidate 12 RWJF Topics into a master categorizer.
  - Test master categorizer against NEW test collection.

# Health Care Coverage > Health Insurance Characteristics sub-queries

(Note: we tested at the whole Topic level, not the sub-queries)

## ❖ Health Care Access

- + Access to Health Care
- + Barriers to Health Care
- + Various Health Outcomes Based on Access

## ❖ Health Care Coverage

- + Health Care Reform – Legislation
- + Health Insurance Characteristics
  - Health Insurance Cost
  - High Deductible Insurance
  - High Risk
  - Mental Health-Substance Abuse
  - PPO-HMO
  - Pre-existing Condition
  - Provider Network-In-Out
- + Health Insurance Enrollment
- + Insurance Marketplace

- ("health care" OR healthcare OR "health insurance") AND (premium\* OR payment\* OR rate\* OR cost\*)
- (high-deductible OR "high deductible") NEAR/5 "health insurance"
- (high-risk OR "high risk") NEAR/5 "health insurance"
- ("mental health" OR "substance abuse") NEAR/5 "health insurance"
- (PPO\* OR HMO\*) NEAR/5 (marketplace\* OR plan\* OR provider\*)
- ("preexisting condition\*" OR "pre-existing condition\*") NEAR/5 "health insurance"
- (("provider network\*" NEAR/5 "health insurance") OR (("in-network" OR "out-of-network") NEAR/5 "health insurance"))

# Optimize recall and precision on a category-by-category basis

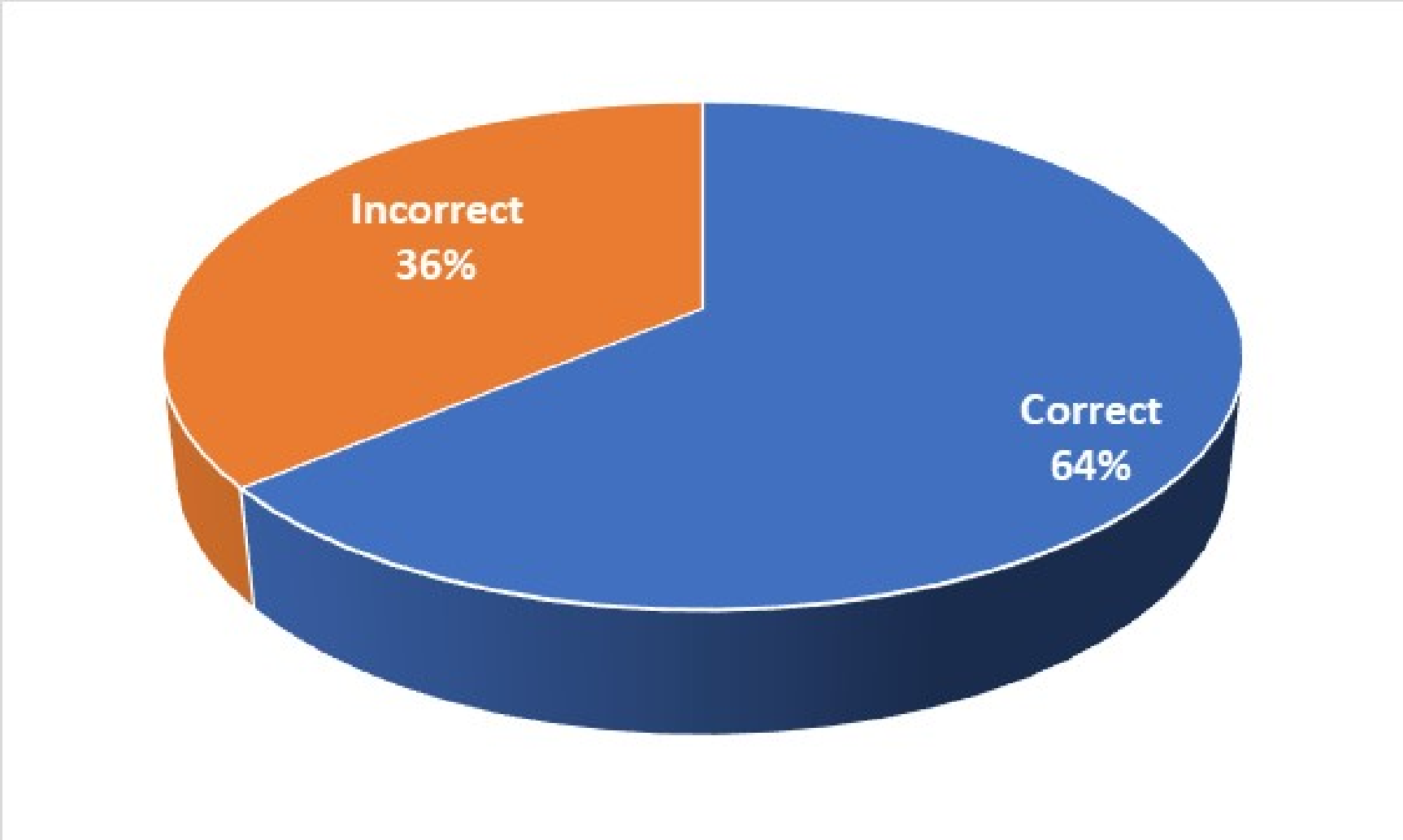
Topic	Recall	Precision
Built Environment & Health	78%	70%
Child & Family Well-being	77%	74%
Childhood Obesity	86%	97%
Disease Prevention & Health Promotion	57%	64%
Early Childhood Development	76%	82%
Health Care Coverage & Access	72%	79%
Health Care Quality & Value	69%	75%
Health Disparities	73%	81%
Health Leadership Development	71%	80%
Nurses and Nursing	94%	94%
Public & Community Health	51%	69%
Social Determinants of Health	71%	69%

## Option 2-Step 3: Consolidate category-by-category queries into a master categorizer

- ❖ Combine all of the category queries into one master categorizer
  - ❖ Make changes to queries as needed to meet the parameters of a single categorizer
  - ❖ Add named entities where appropriate
- ❖ Run analysis and calculate precision and recall to discover if any unintended problems were introduced due to consolidating into one categorizer
  - ❖ Use same test set of assets (Step 2 test set) for Step 3

# Optimize recall and precision for consolidated categorizer

DocId	Core Content	DocId	Build Enrollment	Child Pro	Family Well Being	Obesity	Disease	Prevention	and Health	Promotion	Health Care	Development	Health Care	Quality	and Value	Health	Leadership	Development	Public and Community	Health	Social	Determinants	of Health	Grand Total
1	74742	10	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	12	
0	74955	3	0	0	0	0	0	3	2	2	0	2	0	0	0	0	0	0	0	0	0	0	20	
1	74968	0	2	0	2	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	6	
1	75042	0	0	0	5	0	4	0	41	0	0	0	0	2	0	0	0	0	0	0	0	0	52	
1	75059	1	3	0	4	0	5	2	9	0	3	4	6	37										
1	75083	0	5	0	2	0	23	14	0	0	5	0	0	49										
0	75089	0	2	0	8	0	3	2	0	0	1	0	0	22										
0	75094	0	0	0	0	0	2	0	0	0	0	0	0	2										
0	75095	4	0	0	3	0	10	11	13	0	11	4	0	56										
1	75097	0	0	0	2	0	0	0	0	0	2	0	0	4										
0	75103	0	0	0	0	4	5	0	0	3	0	0	0	21										
0	75109	0	0	0	2	0	1	0	3	0	0	0	0	6										
1	75116	0	5	0	2	5	0	0	0	0	0	0	0	12										
0	75131	1	4	0	2	2	13	3	5	0	3	1	0	34										
1	75148	0	0	0	4	0	3	0	0	0	0	0	0	7										
1	75160	1	0	0	5	1	0	0	3	0	0	0	0	10										
0	75163	2	0	0	0	0	0	6	2	3	0	0	0	23										
1	75173	0	3	0	0	0	0	0	0	1	0	0	0	4										
1	75184	0	0	0	0	0	5	4	0	0	0	0	0	3										
1	75194	0	0	0	3	1	0	0	0	0	0	0	0	4										
1	75196	15	0	0	2	0	0	6	0	4	0	6	33											
1	75198	2	4	0	13	0	15	25	10	2	10	0	81											
0	75199	1	0	0	2	0	4	0	3	0	5	0	18											
1	75201	1	0	0	10	0	7	0	3	3	0	4	28											
1	75202	2	0	0	2	0	6	0	3	3	0	5	21											
0	75209	1	3	0	5	12	3	6	3	3	6	5	59											
0	75213	0	2	0	0	5	0	0	0	0	0	0	7											
0	75214	0	0	0	0	0	0	1	0	0	0	0	1											
1	75215	9	8	5	0	2	0	0	0	0	0	0	4	28										
1	75223	0	0	0	0	0	0	3	0	0	0	0	3											
0	75319	0	0	0	0	0	0	0	2	0	0	0	2											
1	75332	0	2	0	0	0	0	0	0	0	0	4	6											
0	75333	0	5	0	0	0	0	0	0	0	0	0	5											
1	75362	8	0	0	0	0	0	3	0	4	0	0	15											
1	75373	0	0	0	3	0	4	9	0	0	0	0	16											
1	75374	0	0	0	0	0	7	0	0	0	0	0	7											
0	75386	2	3	0	0	0	0	3	3	0	0	2	13											
1	75403	0	0	0	0	0	21	14	0	0	0	0	37											
1	75417	0	11	0	2	3	0	0	0	0	5	0	21											
1	75432	7	6	4	0	2	0	0	0	0	0	0	23											
0	75449	0	6	0	5	8	0	0	8	0	0	0	32											
1	75475	0	0	0	0	5	0	0	0	0	0	0	5											
1	75493	0	2	0	4	2	0	1	3	0	3	4	19											
1	75520	2	0	0	2	0	0	0	2	0	0	0	15											
1	75523	0	0	0	0	0	0	0	3	0	4	0	7											
1	75531	11	0	0	3	0	3	4	0	0	0	0	32											
0	75541	2	0	0	0	0	2	0	3	0	0	5	12											



... is this performance reliable enough for unsupervised categorization of high risk content?

## Option 3: Combination of mapping and automated methods

- ❖ Used mapping to provide initial conversion of legacy content from A-Z to new faceted taxonomy.
- ❖ Identified 1) categories with no mapping to any grants, and 2) grants with too many Topics; then processed these grants with a combination of human review, and automated categorization to improve their indexing.
- ❖ Worked to build a critical mass of new grants indexed by editors using the new taxonomy.
- ❖ **Future goal:** Automation-assisted indexing.

# Summary

- ❖ Natural language processing enables unsupervised automated categorization.
- ❖ Individual terms can be replaced by pre-defined queries chosen and constructed by editors.
- ❖ Good results on individual pre-defined queries/topics may not be as good when distinguishing which is most relevant among a group of pre-defined queries/topics.
- ❖ Automated categorization results may be good enough for some collections of content (high volume or low risk), but not for others (low value or high risk).



# Questions

Joseph Busch

[jbusch@taxonomystrategies.com](mailto:jbusch@taxonomystrategies.com) or [joseph@semanticstaffing.com](mailto:joseph@semanticstaffing.com)

@joebusch

mobile +1 415-377-7912

# Supervised vs. Unsupervised Automated Categorization: Summary

- ❖ AI promises to categorize all types of content with reliable results, but the reality is much more complex. Most applications won't work with a meat grinder approach where you pour a huge amount of content in one end and a perfectly organized collection comes out the other end. Effective automated categorization depends on defining a process workflow and assembling a stack of methods to process different types of content in different ways. Designing and validating a content processing workflow requires human judgements. So good quality categorization applications often rely on how to make the best use of people. This presentation provides a reality check on unsupervised automated categorization, and discusses a case study where the performance was suitable for editorial review and approval, but not for unsupervised processing of a large collection. The project is the Robert Wood Johnson Foundation grants information system using Lexalytics configurations developed by Taxonomy Strategies.