



Systematic Characterization of the NIH Prevention Research Portfolio

Patricia Mabry¹, Jennifer Villani¹, Jessica Wu^{1*}, Jocelyn Lee¹, Ranell Myles¹, Pamela Carter-Nolan², Lamyaa Yousif², Richard Panzer², Jason Hamrick², Ifeyinwa Udo², David M. Murray¹

¹Office of Disease Prevention, Division of Program Coordination, Planning, and Strategic Initiatives, National Institutes of Health and ²IQ Solutions, Inc., Rockville, MD

*Corresponding author: Jessica.wu2@nih.gov

Goal

The first strategic priority of the NIH Office of Disease Prevention (ODP) Strategic Plan for 2014-2018 is to **systematically monitor NIH investments in prevention research and assess the progress and results of that research**. The ODP is developing a computer-based automated tool capable of objectively classifying the NIH prevention research grant portfolio on a number of dimensions of interest.

Project Rationale and Summary

Automated portfolio analysis tools at NIH include the RCDC System, which categorizes and reports NIH spending in 233 categories. While Prevention is one of those categories, further characterization of grants within the Prevention category is limited. The ODP has identified approximately 150 study characteristics (topics) across eight categories that are of interest to the NIH prevention research community. Manual categorization is impractical on a large scale and human factors inherent in manual coding lead to inconsistencies in coding across coders and time. Therefore, the ODP seeks to develop a computer-based solution to enable an automated, standardized, rapid, and objective characterization of NIH prevention research funding. This tool will enable identification of patterns and trends in NIH prevention research funding and research areas that may benefit from targeted investments by the NIH Institutes and Centers. Such categorization will enable assessment of the progress and changes in NIH-funded prevention research over time to inform program planning and reporting.

Developing this tool is a multi-step process, described here. Four partner organizations are participating – the NIH ODP; the ODP's contractor, IQ Solutions, Inc. (IQS); the NIH Office of Portfolio Analysis (OPA); and the NIH Center for Information Technology (CIT). The development of the initial tool is expected to take at least two years to complete and involves:

- ◆ Developing a taxonomy, or framework, for coding abstracts
- ◆ Manually coding grant abstracts according to the taxonomy to develop a “training set”
- ◆ Using a machine learning tool which uses a “training set” of abstracts to generate an algorithm for automated classification, the **Portfolio Learning Tool (PLT)**
- ◆ Using the PLT to code new abstracts and manually validating that output

The Taxonomy

The ODP has developed a prevention research taxonomy for coding grant abstracts along with a detailed protocol. The protocol ensures consistent application of the Prevention Taxonomy across abstracts and coders.

- ◆ **8 taxonomy categories:** Study focus (composed of Study Rationale, Exposure variables, and Outcome variables), Entities studied, Study setting, Population focus, Study design, and Prevention research category
- ◆ Within each category are **topics**; one or more topics may apply to any given abstract
- ◆ Taxonomy **protocol:** provides detailed instructions, definitions, and examples to facilitate correct, standardized use of the Prevention Taxonomy by coders.

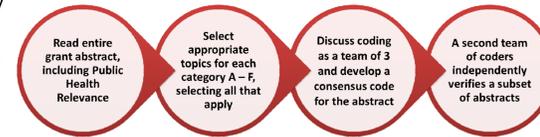
Prevention Research Taxonomy

Team Coding of Grant Abstracts

Generating the training set for the PLT is the most time consuming and staff intensive aspect of the project. Three coders simultaneously read a funded NIH grant abstract and independently select applicable topics for each of the eight categories on the taxonomy form, following protocol instructions. The group discusses their selections, and develops a fourth “consensus” set of selections. A second team of three independently codes a subset of the abstracts to ensure strict compliance with the protocol.

To develop the training set, thousands of NIH Type 1 R01 grant abstracts will be coded. The **Prevention Abstract Classification Tool (PACT)** is custom software developed by ODP and IQS to facilitate the team coding process and data collection. PACT features include:

- ◆ Record/track assignment of coders to teams (team membership is reconfigured daily)
- ◆ Assign abstracts to teams for coding and flag abstracts for verification coding
- ◆ Capture coding selections on taxonomy form using tablets such as iPads
- ◆ Calculate inter-rater reliability data by coder, abstract, and team across each of eight categories—104 Kappa statistics per abstract
- ◆ Provide searchable taxonomy protocol
- ◆ Compare individual and team performance to established criteria
- ◆ Export data and generate reports, including information on progress and productivity



Coding a training set of exemplars using PACT

1. First, an IQS team of three individually codes an abstract, referencing the protocol when necessary, and enters their codes into PACT.
2. Once each member of the team has submitted their codes in PACT, the team discusses their coding and reconciles any differences to generate a consensus set of codes.
3. To verify the accuracy of the coding, a team of 3 coders from the ODP codes the same abstract, using the same process as the IQS team.
4. After the ODP team has individually coded an abstract and developed a consensus, it also compares the ODP consensus codes against the IQS consensus codes, and then reconciles any differences to generate a “final consensus.” The ODP team codes a subset of abstracts completed by IQS.
5. The PACT system records the codes for both individual and consensus codes. The “final consensus” codes will be used as exemplars in the training set for the PLT.

Individual Coding

Consensus Coding in Progress

Completed Consensus Coding

Top: Individual coding, Middle: Consensus coding in progress, Bottom: Consensus coding. Color coding: dark green indicates consensus selections, bright green indicates all 3 coders made the selection, yellow indicates 2 coders made that selection, and red indicates only 1 coder made that selection.

Coder Training and Certification

Accurate coding is an acquired skill requiring a good deal of practice across a wide range of research topics. Therefore, all coders must complete an extensive training program and meet individual performance criteria on a standard set of abstracts prior to being certified to generate exemplars. The training/certification process includes:

- ◆ Training materials: syllabus, annotated abstracts, slide shows, worksheets
- ◆ Videos demonstrating consensus discussions
- ◆ Over 150 practice abstracts with answer keys, reviewed with certified coders
- ◆ Practice developing consensus with certified coders
- ◆ Standard sets of abstracts for self-assessment and certification

Coders typically train in cohorts of 2-6 people. An integral part of the process includes ample discussion and review of keys to ensure that each coder understands specific abstract text leading to coding decisions and correct interpretation of protocol language.

Inter-Rater Reliability and Coder Certification

PACT is also integrated with SAS to calculate inter-rater reliability between each coder and the final consensus for each category for every abstract. Specialized Kappa statistics are used to evaluate coder and team performance and for identifying difficult-to-code abstracts. Category-specific performance (Kappa values) by ODP staff trained in the use of the Prevention Taxonomy and protocol were used as the basis for computing threshold values for the Kappa statistics. These thresholds serve as the criteria for coder certification. Kappas also aid in coder self-assessment and in continuous monitoring of both individual and team performance.

Complete set of Kappa statistics calculated for one abstract coded by both IQS and ODP teams.

Progress to Date and Next Steps

Ultimately, this work will produce an automated tool designed to facilitate the identification of patterns and trends in NIH prevention research funding and research areas that may benefit from targeted investments by the NIH Institutes and Centers.

Progress to date:

- ◆ The Prevention Taxonomy and protocol have been developed and refined
- ◆ A team coding approach and quality control process have been developed and refined
- ◆ An extensive training program and certification program have been developed and piloted
- ◆ Threshold Kappa values have been established to certify coders and assess performance
- ◆ PACT software has been developed, tested, and implemented; enhancements are ongoing
- ◆ A total of 1100 abstracts have been coded to date

Next steps:

- ◆ Certify additional trainees and continue to grow the training set for the PLT
- ◆ Share all relevant materials and processes with interested parties, particularly prevention-focused NIH staff. This includes the Prevention Taxonomy and protocol, team coding approach, training and certification materials, and PACT software.

Acknowledgements

ODP: Jody Engel, Denise Simons-Morton

IQS: Kimberly James, Shahina Akter, Jimmy Ngo, Anthony Ortiz, Daniela Poss, Agnieszka Roman, Alexis Hall, Arielle Dolegui, Adeola Olufunmilade, Ifeyinwa Udo, Lindsey Beattie, Kossia Dassie, Teemar Fisseha, Lolita Kachay, Eric Nwazue, Anne Abbate, Andrea Kelsey

IQS PACT team: Richard Panzer, Jason Hamrick, Jasmine Chan, Lucy Liu, Tatiana Timonina, Jason Summers, Sergey Kukin, Sunil Taneja

Prevention Portfolio Analysis Tool Development with the Portfolio Learning Tool

Develop a taxonomy for coding abstracts

Code a training set of exemplars

Use manually coded data to inform the PLT

Validate the PLT output

Goal: PLT can accurately retrieve grants according to taxonomy