

1 **Design and methodological considerations for biomarker discovery**
2 **and validation in the Integrative Analysis of Lung Cancer Etiology**
3 **and Risk (INTEGRAL) Program**
4

5 Hilary A Robbins PhD¹, Karine Alcalá MS^{^1}, Elham Khodayari Moez PhD^{^2}, Florence Guida
6 PhD³, Sera Thomas MSc², Hana Zahed MS¹, Matthew T Warkentin MSc^{2,4}, Karl Smith-Byrne
7 DPhil⁵, Yonathan Brhane MS^{2,4}, David Muller PhD⁶, Xiaoshuang Feng PhD¹, Demetrius
8 Albanes MD⁷, Melinda C Aldrich PhD⁸, Alan A Arslan MD⁹, Julie Bassett PhD¹⁰, Christine D
9 Berg MD¹¹, Qiuyin Cai MD PhD¹², Chu Chen PhD¹³, Michael PA Davies PhD¹⁴, Brenda
10 Diergaarde PhD^{15,16}, John K Field PhD¹⁴, Neal D Freedman PhD⁷, Wen-Yi Huang PhD⁷,
11 Mikael Johansson MD¹⁷, Michael Jones PhD¹⁸, Woon-Puay Koh MBBS PhD^{19,20}, Stephen
12 Lam MD²¹, Qing Lan MD PhD⁷, Arnulf Langhammer MD PhD^{22,23}, Linda M Liao PhD⁷,
13 Geoffrey Liu MD²⁴, Reza Malekzadeh MD²⁵, Roger L Milne PhD^{10,26,27}, Luis M Montuenga
14 PhD^{28,29,30}, Thomas Rohan MBBS PhD³¹, Howard D Sesso ScD³², Gianluca Severi PhD³³,
15 Mahdi Sheikh MD PhD¹, Rashmi Sinha PhD⁷, Xiao-Ou Shu MD PhD¹², Victoria L Stevens
16 PhD³⁴, Martin C Tammemägi DVM PhD^{35,36}, Lesley F Tinker PhD³⁷, Kala Visvanathan MD
17 MHS³⁸, Ying Wang PhD³⁹, Renwei Wang MD⁴⁰, Stephanie J Weinstein PhD⁷, Emily White
18 PhD⁴¹, David Wilson MD MPH⁴², Jian-Min Yuan MD PhD^{43,16}, Xuehong Zhang PhD³², Wei
19 Zheng MD PhD¹², Christopher I Amos PhD⁴⁴, Paul Brennan PhD¹, Mattias Johansson PhD^{*1},
20 Rayjean J Hung PhD^{*2,4}

21
22 [^]Contributed equally (KA, EKM)

23 ^{*}Joint senior authors (MJ, RJH)

24
25 ¹Genomic Epidemiology Branch, International Agency for Research on Cancer, Lyon, France,
26 ²Prosserman Centre for Population Health Research, Lunenfeld-Tanenbaum Research Institute, Sinai
27 Health, Toronto, Canada, ³Environment and Lifestyle Epidemiology Branch, International Agency for
28 Research on Cancer, Lyon, France, ⁴Dalla Lana School of Public Health, University of Toronto,
29 Toronto, Canada, ⁵Cancer Epidemiology Unit, University of Oxford, Oxford, United Kingdom, ⁶Division
30 of Genetic Medicine, Imperial College London School of Public Health, London, United Kingdom,
31 ⁷Division of Cancer Epidemiology and Genetics, National Cancer Institute, Rockville, MD, USA,
32 ⁸Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA, ⁹Departments of
33 Obstetrics and Gynecology and Population Health, New York University Grossman School of
34 Medicine, New York, NY, USA, ¹⁰Cancer Epidemiology Division, Cancer Council Victoria, Melbourne,
35 Australia, ¹¹Retired, Bethesda, MD, USA, ¹²Division of Epidemiology, Department of Medicine,
36 Vanderbilt University Medical Center, Nashville, TN, USA, ¹³Program in Epidemiology and the
37 Women's Health Initiative Clinical Coordinating Center, Division of Public Health Sciences, Fred

38 Hutchinson Cancer Research Center, Seattle, WA, USA, ¹⁴Molecular & Clinical Cancer Medicine,
39 University of Liverpool, Liverpool, United Kingdom, ¹⁵Department of Human Genetics, Graduate
40 School of Public Health, University of Pittsburgh, Pittsburgh, PA, USA, ¹⁶UPMC Hillman Cancer
41 Centre, Pittsburgh, PA, USA, ¹⁷Department of Radiation Sciences, Oncology, Umea University,
42 Umea, Sweden, ¹⁸Division of Genetics and Epidemiology, Institute of Cancer Research, London,
43 United Kingdom, ¹⁹Healthy Longevity Translational Research Program, Yong Loo Lin School of
44 Medicine, National University of Singapore, Singapore, Singapore, ²⁰Singapore Institute for Clinical
45 Sciences, Agency for Science Technology and Research (A*STAR), Singapore, Singapore,
46 ²¹Integrative Oncology, British Columbia Cancer Agency, Vancouver, Canada, ²²HUNT Research
47 Center, Department of Public Health and Nursing, NTNU Norwegian University of Science and
48 Technology, Levanger, Norway, ²³Levanger Hospital, Nord-Trøndelag Hospital Trust, Levanger,
49 Norway, ²⁴Computational Biology and Medicine Program, Princess Margaret Cancer Center, Toronto,
50 Canada, ²⁵Digestive Disease Research Institute, Tehran University of Medical Sciences, Tehran, Iran,
51 ²⁶Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health,
52 University of Melbourne, Parkville, Australia, ²⁷School of Clinical Sciences at Monash Health, Monash
53 University, Melbourne, Australia, ²⁸Center of Applied Medical Research (CIMA) and Schools of
54 Sciences and Medicine, University of Navarra, Pamplona, Spain, ²⁹IDISNA, Pamplona, Spain,
55 ³⁰CIBERONC, Madrid, Spain, ³¹Department of Epidemiology & Population Health, Albert Einstein
56 College of Medicine, Bronx, NY, USA, ³²Brigham and Women's Hospital, Harvard Medical School,
57 Boston, MA, USA, ³³Inserm, Université Paris-Saclay, Villejuif, France, ³⁴Rollins School of Public
58 Health, Emory University, Atlanta, GA, USA, ³⁵Department of Health Sciences, Brock University, St.
59 Catharines, ON, Canada, ³⁶Prevention and Cancer Control, Ontario Health, Toronto, ON, Canada,
60 ³⁷Women's Health Initiative Clinical Coordinating Center, Division of Public Health Sciences, Fred
61 Hutchinson Cancer Research Center, Seattle, WA, USA, ³⁸Department of Epidemiology, Johns
62 Hopkins Bloomberg School of Public Health, Baltimore, MD, USA, ³⁹American Cancer Society,
63 Atlanta, GA, USA, ⁴⁰UPMC Hillman Cancer Center, University of Pittsburgh, Pittsburgh, PA, USA,
64 ⁴¹Cancer Prevention Research Program, Fred Hutchinson Cancer Research Center, Seattle, WA,
65 USA, ⁴²Division of Pulmonary, Allergy, and Critical Care Medicine, Department of Medicine, University
66 of Pittsburgh, Pittsburgh, PA, USA, ⁴³Department of Epidemiology, Graduate School of Public Health,
67 University of Pittsburgh, Pittsburgh, PA, USA, ⁴⁴Institute for Clinical and Translational Research,
68 Baylor College of Medicine, Houston, TX, USA

69

70 **Corresponding authors:**

71 Hilary Robbins (RobbinsH@iarc.fr) and Mattias Johansson (JohanssonM@iarc.fr)

72 Genomic Epidemiology Branch

73 International Agency for Research on Cancer

74 150 cours Albert Thomas

75 CEDEX 69732 Lyon, France

76

77 Rayjean Hung (Rayjean.hung@lunenfeld.ca)

78 Lunenfeld-Tanenbaum Research Institute, Sinai Health

79 Dalla Lana School of Public Health, University of Toronto,

80 60 Murray St. Toronto, ON M5T 3L9. Canada

81

82 **Conflicts of interest:**

83 Dr Montuenga reports the following potential conflicts of interest: Astra-Zeneca (speaker's
84 bureau and research grant), Bristol Myers Squibb (research grant), AMADIX: (licensed
85 patent co-holder on complement fragments for lung cancer early detection).

86 All other authors report no conflicts of interest.

87

88 **Funding:**

89 This study was supported by the US NCI (INTEGRAL program U19 CA203654 and R03
90 CA245979), the Lung Cancer Research Foundation, l'Institut National Du Cancer (2019-1-
91 TABAC-01, INCa, France), the Cancer Research Foundation of Northern Sweden (AMP19-
92 962), and an early detection of cancer development grant from Swedish Department of
93 Health ministry. RJH is supported by the Canada Research Chair of the Canadian Institute
94 of Health Research. LMM was supported by FIMA, Fundación ARECES, ISCIII-Fondo de
95 Investigación Sanitaria-Fondo Europeo de Desarrollo Regional (PI19/00098) and a grant
96 from The Lung Ambition Alliance. MCA is supported by NCI R01 CA251758. The ATBC
97 Study is supported by the Intramural Research Program of the U.S. National Cancer
98 Institute, National Institutes of Health, Department of Health and Human Services. The
99 Southern Community Cohort Study was supported by NCI U01CA202979. The Physicians'
100 Health Study (PHS) is supported by research grants CA097193, CA34944, CA40360,
101 HL26490, and HL34595 from the NIH. The Women's Health Study (WHS) is supported by
102 research grants EY06633, EY18820, CA047988, HL043851, HL080467, HL099355, and
103 CA182913 from the NIH. The WHI program is funded by the National Heart, Lung, and Blood
104 Institute, National Institutes of Health, U.S. Department of Health and Human Services
105 through contracts 75N92021D00001, 75N92021D00002, 75N92021D00003,
106 75N92021D00004, 75N92021D00005. CLUE II funding was from the National Cancer
107 Institute (U01 CA86308, Early Detection Research Network; P30 CA006973), National
108 Institute on Aging (U01 AG18033), and the American Institute for Cancer Research.
109 Maryland Cancer Registry (MCR) Cancer data was provided by the Maryland Cancer
110 Registry, Center for Cancer Prevention and Control, Maryland Department of Health, with

111 funding from the State of Maryland and the Maryland Cigarette Restitution Fund. The
112 collection and availability of cancer registry data is also supported by the Cooperative
113 Agreement NU58DP006333, funded by the Centers for Disease Control and Prevention.
114 Acknowledgements for the NIH-AARP study are available at:
115 <https://dietandhealth.cancer.gov/acknowledgement.html>. P LuSS was supported by NCI P50
116 CA090440 and NCI P30 CA047904.

117

118 **Data sharing statement:**

119 Researchers who are interested in analyzing the Lung Cancer Cohort Consortium (LC3)
120 dataset are encouraged to contact Dr. Robbins or Dr. Johansson. The LC3 Access Policy is
121 available at the following link: [https://www.iarc.who.int/wp-](https://www.iarc.who.int/wp-content/uploads/2021/12/LC3_Access_Policy.pdf)
122 [content/uploads/2021/12/LC3_Access_Policy.pdf](https://www.iarc.who.int/wp-content/uploads/2021/12/LC3_Access_Policy.pdf).

123

124 **Disclaimer:**

125 Where authors are identified as personnel of the International Agency for Research on
126 Cancer / World Health Organization, the authors alone are responsible for the views
127 expressed in this article and they do not necessarily represent the decisions, policy, or views
128 of the International Agency for Research on Cancer / World Health Organization. The
129 contents of this manuscript are solely the responsibility of the authors and do not necessarily
130 represent the official views of the Centers for Disease Control and Prevention or the
131 Department of Health and Human Services, nor does mention of trade names, commercial
132 products, or organizations imply endorsement by the US government.

133

134 **Keywords:**

135 Lung cancer screening, early detection, biomarkers, risk prediction, nodule malignancy,
136 biomarker discovery and validation, study design

137

138 **Word counts:**

139 Abstract: 250

140 Text: 3958

141 Tables/figures: 5

142 Supplement: 2 tables and additional text

143

144 **Design and methodological considerations for biomarker discovery**
145 **and validation in the Integrative Analysis of Lung Cancer Etiology**
146 **and Risk (INTEGRAL) Program**
147

148

149

Keywords:

150 Lung cancer screening, early detection, biomarkers, risk prediction, nodule malignancy,

151 biomarker discovery and validation, study design

152

Word counts:

154 Abstract: 250

155 Text: 3958

156 Tables/figures: 5

157 Supplement: 2 tables and additional text

158 **Abstract**

159

160 The Integrative Analysis of Lung Cancer Etiology and Risk (INTEGRAL) program is an NCI-
161 funded initiative with an objective to develop tools to optimize lung cancer screening. Here,
162 we describe the rationale and design for the Risk Biomarker and Nodule Malignancy projects
163 within INTEGRAL.

164

165 The overarching goal of these projects is to systematically investigate circulating protein
166 markers to include on a panel for use (i) pre-LDCT, to identify people likely to benefit from
167 screening, and (ii) post-LDCT, to differentiate benign versus malignant nodules. To identify
168 informative proteins, the Risk Biomarker project measured 1,161 proteins in a nested-case
169 control study within 2 prospective cohorts (n=252 lung cancer cases and 252 controls) and
170 replicated associations for a subset of proteins in 4 cohorts (n=479 cases and 479 controls).
171 Eligible participants had any history of smoking and cases were diagnosed up to 3 years
172 following blood draw. The Nodule Malignancy project measured 1,078 proteins among
173 participants with a heavy smoking history within 4 LDCT screening studies (n=425 cases
174 diagnosed up to 5 years following blood draw, 398 benign-nodule controls, and 430 nodule-
175 free controls).

176

177 The INTEGRAL panel will enable absolute quantification of 21 proteins. We will evaluate its
178 performance in the Risk Biomarker project using a case-cohort study including 14 cohorts
179 (n=1,696 cases and 2,926 subcohort representatives), and in the Nodule Malignancy project
180 within 5 LDCT screening studies (n=675 cases, 648 benign-nodule controls, and 680
181 nodule-free controls). Future progress to advance lung cancer early detection biomarkers
182 will require carefully designed validation, translational, and comparative studies.

183 **Introduction**

184

185 Lung cancer screening by low-dose computed tomography (LDCT) has accelerated the field
186 of lung cancer research with a renewed focus on early detection.^{1,2} However, several
187 questions remain regarding how to best implement LDCT screening,³ including how to
188 identify individuals who are likely to benefit from screening, and how to manage nodules of
189 indeterminate malignancy status identified on LDCT scans. Here, we describe the rationale
190 and design of a large international research effort to develop and validate biomarker tools
191 that can be applied in these two settings.

192

193 In 2018, the US National Cancer Institute (NCI) funded the Integrative Analysis of Cancer
194 Risk and Etiology (INTEGRAL) U19 program, which includes an objective to develop early
195 detection biomarkers and risk prediction tools for lung cancer screening. The INTEGRAL
196 program comprises 3 projects: the Genetics project which studies germline genetics, the
197 Risk Biomarker project which studies pre-diagnostic blood biomarkers, and the Nodule
198 Malignancy project which studies applications in LDCT screening studies including nodule
199 evaluation. Here, we describe a joint effort of the Risk Biomarker and Nodule Malignancy
200 projects to systematically investigate circulating protein markers for pre- and post-LDCT
201 applications.

202

203 The primary objective of the Risk Biomarker project is to identify and validate biomarkers
204 that can improve lung cancer risk prediction among people with a smoking history. A
205 secondary objective is to develop and validate questionnaire-based lung cancer risk
206 prediction models. The objectives for the Nodule Malignancy project are to identify
207 biomarkers and establish quantitative imaging models that can differentiate benign versus
208 malignant nodules following an initial LDCT scan. The Risk Biomarker project leverages
209 resources from the Lung Cancer Cohort Consortium (LC3)⁴⁻⁸ which was initially established
210 in 2010 within the NCI Cohort Consortium.⁹ The Nodule Malignancy project brings together
211 LDCT screening studies in the framework of the International Lung Cancer Consortium
212 (ILCCO), which has provided a foundation for collaborative research on lung cancer since
213 2004 (<http://ilcco.iarc.fr>).

214

215 This paper provides a design overview of the biomarker studies within the INTEGRAL Risk
216 Biomarker and Nodule Malignancy projects. We highlight considerations that motivated the

217 design, present details of the study population, and describe the harmonized databases
218 resulting from these projects. Finally, we discuss perspectives for research to follow this
219 initiative with a view toward implementation of the prediction tools in clinical practice.

220

221 **Development and validation of a protein biomarker panel for early lung** 222 **cancer detection**

223

224 **Motivation**

225

226 The US Preventive Services Task Force (USPSTF) currently recommends lung cancer
227 screening for people aged 50-80 years who have smoked at least 20 pack-years and
228 currently smoke or have quit within the past 15 years.¹⁰ However, more than one-third of
229 lung cancer deaths that could be prevented among people who have smoked fall outside of
230 these criteria.¹¹ To better target the highest-risk population, screening can instead be offered
231 to people whose individual lung cancer risk exceeds a certain threshold as estimated by a
232 risk prediction model.¹²⁻¹⁵ This approach is included in the US National Comprehensive
233 Cancer Network (NCCN) guidelines.¹⁶

234

235 Biomarkers may provide additional or complementary information on lung cancer risk and
236 represent a promising avenue to improve existing risk prediction models. Conceptually, this
237 could improve efficiency in two ways: by offering screening to people who have high risk
238 based on biomarkers but are not otherwise eligible for screening based on the current
239 recommendation, and by deprioritizing screening for individuals who are eligible but have a
240 low-risk biomarker profile. Various domains of biomarkers have been investigated, but the
241 translation of this research into practice has been slow, partly due to the lack of appropriately
242 designed studies to establish and validate biomarker-based risk prediction models.^{17,18}

243

244 Another setting in which biomarkers could be applied in lung cancer screening is to better
245 distinguish between malignant and benign nodules on LDCT images. Nodules are detected
246 in up to one-quarter of participants, but the vast majority are benign. Managing nodules with
247 uncertain clinical significance (i.e., indeterminate nodules) represents an important challenge
248 because false-positive nodules can lead to interventions with risks of long-term harm. On the
249 other hand, missed malignant nodules can lead to a lost opportunity for curative treatment.

250 Several prediction models for nodule malignancy have been developed,^{19–21} but their
251 classification accuracies remain imperfect.

252

253 Recent papers have highlighted common limitations in the design of studies aiming to
254 identify and validate biomarkers for early cancer detection,²² including lung cancer.¹⁸ To
255 avoid common biases resulting from systematic differences between cases and controls, the
256 prospective-specimen-collection, retrospective-blinded-evaluation (PRoBE) design
257 emphasizes the use of pre-diagnostic samples, sampling from the same source population,
258 and matching on important factors that impact biomarker measurements and outcome.²³ In
259 validation studies, it is critical that the added contribution of the biomarker, compared with
260 existing tools, can be clearly identified and quantified.¹⁸

261

262 Several studies led by our group and others informed our overall choice to pursue a
263 research program focused on protein biomarkers within INTEGRAL. First, in a pilot study
264 published in 2018, members of our team found that a pre-defined set of cancer-related
265 protein biomarkers improved discrimination between lung cancer cases and controls
266 compared to a smoking-based risk prediction model, when the markers were measured in an
267 independent validation study using samples collected within the year before diagnosis.²⁴
268 Second, we carried out a modeling study which suggested that using such biomarkers to
269 optimize screening eligibility could be cost-effective, as long as the biomarker provides
270 moderate or better risk discrimination at modest cost.²⁵ Studies also suggest that protein
271 markers can improve discrimination between malignant and benign lung nodules.^{26,27}
272 Therefore, building on these promising preliminary data, the INTEGRAL program was
273 formed to conduct a comprehensive protein biomarker evaluation from discovery to
274 validation for both population-based risk prediction (Risk Biomarker project), and nodule
275 differentiation (Nodule Malignancy project).

276

277 Our overarching aims are *i)* to identify circulating proteins that provide additional information
278 to the gold standard on both lung cancer risk and nodule malignancy and *ii)* to develop and
279 validate a multiplex lung cancer biomarker assay that can quantify key lung cancer risk
280 and/or nodule malignancy proteins in small volumes of peripheral blood in a cost-effective
281 manner. Use of a single assay will help to streamline clinical implementation along the
282 various steps of the LDCT screening pathway.

283

284 **Design**

285

286 Overview

287 **Figure 1** outlines the sequential study phases of the INTEGRAL Risk Biomarker and Nodule
288 Malignancy projects. In the Risk Biomarker project, using pre-diagnostic samples from
289 population cohorts, an initial ‘full discovery’ phase scanned a broad set of protein markers,
290 followed by a ‘targeted discovery’ phase which replicated results for a subset of proteins.
291 The Nodule Malignancy project started with an expanded targeted discovery phase and
292 analyzed samples from LDCT screening studies to identify proteins that are specifically
293 useful to distinguish between benign and malignant lung nodules. The results from both
294 projects will be used to configure the INTEGRAL panel with 21 circulating protein markers,
295 whose performance will be assessed in a validation phase conducted separately within each
296 project. **Table 1** summarizes the key characteristics of the participating cohorts and LDCT
297 screening studies in each phase.

298

299 We are using the Olink proteomics platform (Olink Proteomics, Uppsala, Sweden)
300 throughout the project.²⁸ Olink discovery assays allow high-throughput semi-quantified
301 concentration measures of highly annotated proteins in less than 50 uL of plasma or serum.
302 The platform uses proximity extension assay (PEA) technology which is highly sensitive,
303 avoids cross-reactivity, and has high reproducibility. Relative protein concentrations are
304 expressed as normalized protein expression (NPX) on log₂ scale, which is estimated from
305 quantitative PCR cycle threshold values, and were standardized for analysis. For all
306 laboratory analyses in INTEGRAL, cases and controls are randomly allocated across plates,
307 with matched pairs plated together where relevant.

308

309 To enable absolute quantification of proteins for clinical applications, we will develop the
310 INTEGRAL panel as an Olink customized panel. Customized panels are also based on PEA
311 technology and can measure up to 21 proteins in less than 50 uL of plasma or serum.²⁹ We
312 plan to include 21 proteins on our panel, which is the maximum due to technical limitations,
313 since reducing the number of proteins reduces neither the assay cost nor the sample volume
314 requirement.

315

316 Risk Biomarker project

317 The design of the Risk Biomarker project was informed by several considerations. Given that
318 a key application for biomarkers in screening eligibility could be to identify individuals at high

319 risk for lung cancer despite not meeting eligibility criteria (e.g., USPSTF criteria), it was
320 crucial that the Risk Biomarker project include individuals who are both eligible and ineligible
321 by current criteria. Therefore, pre-diagnostic samples collected within prospective cohorts
322 provided an ideal study resource. Within cohorts, we first restricted to participants who
323 currently or formerly smoked because they represent the current target population for lung
324 cancer screening.¹⁰ Second, we included cases diagnosed up to 3 years following blood
325 draw, to predict lung cancer within a clinically actionable timeframe.²⁴ Third, we used a
326 matched case-control design for the discovery phases, but a case-cohort design for the
327 validation phase. For discovery, the matched design is important to eliminate influences
328 such as storage duration and biospecimen handling. In the validation phase, we changed to
329 a case-cohort design to facilitate development of an integrated risk prediction model that is
330 well-calibrated and representative of the source population (i.e., representative of all
331 participants in the cohorts who ever smoked).

332

333 *Full discovery phase*

334 In the Risk Biomarker project full discovery phase, we measured all 13 Olink proteomics
335 panels available in late 2019, which cover a range of domains including inflammation,
336 oncology, and cardiovascular disease (1,161 proteins, **Appendix Table, Table 2**). The
337 objective of the full discovery phase was to select panels to measure in the targeted
338 discovery phase, and the sample included the European Investigation into Cancer and
339 Nutrition (EPIC, n=188 lung cancer cases) and the Northern Sweden Health and Disease
340 Study (NSHDS, n=64 cases) (**Table 1**; further details in **Supplementary Table 1**). We
341 included all confirmed lung cancer cases among people who ever smoked that were
342 diagnosed within 3 years of blood draw. For each case, one control was randomly chosen
343 using incidence density sampling from risk sets consisting of people who ever smoked and
344 were alive and free of cancer at the time of diagnosis of the index case. Matching criteria
345 included cohort, study center (where relevant), sex, date of blood collection (within 1 month
346 of the index case, relaxed to 3 months for cases without available controls), date of birth
347 (within 1 year of the index case, relaxed to 3 years), and smoking status in 4 categories:
348 people who formerly smoked and quit <10 or ≥10 years prior, and people who currently
349 smoked <15 or ≥15 cigarettes per day.

350

351 The dataset generated by the full discovery phase therefore includes 252 case-control pairs
352 with 1,161 proteins measured on each participant (**Table 2**). Statistical analyses applied

353 conditional logistic and penalized regression. We used the results to examine, for each of
354 the 13 proteomics panels, the number of highly ranked and consistently selected proteins.

355

356 *Targeted discovery phase*

357 The targeted discovery phase of the Risk Biomarker project used the same design to
358 independently replicate associations for a subset of proteomics panels, chosen to maximize
359 coverage of the promising proteins while minimizing the total cost. This phase included 4
360 cohorts with 479 total eligible lung cancer cases: the Cancer Prevention Study II, the Nord-
361 Trøndelag Health Study, the Melbourne Collaborative Cohort Study, and the Singapore
362 Chinese Health Study (**Table 1**; further details in **Supplementary Table 1**). To cover as
363 many of the promising proteins as possible, we measured the Immuno-oncology,
364 Oncology II, Cardiovascular III, and Inflammation panels on all 4 cohorts, and the
365 Oncology III and Neuro-exploratory panels on 3 cohorts each (**Table 2**).

366

367 The dataset generated for the targeted discovery phase therefore includes 479 case-control
368 pairs with between 392 and 484 proteins measured for each participant (**Table 2**). Statistical
369 analyses included conditional logistic regression, penalized regression, and stratified
370 approaches. For the INTEGRAL panel, we are prioritizing proteins selected in penalized
371 regression models that show a consistent association with lung cancer across cohorts.

372

373 *Validation phase*

374 The Risk Biomarker project validation phase includes 14 cohorts and employs a case-cohort
375 design. In each cohort, all cases diagnosed within 3 years of blood draw were included.
376 Subcohort representatives were randomly sampled at the time of blood draw in 8 jointly
377 defined categories including age (above or below the median age among cases), sex (male
378 or female, except for single-sex cohorts), and smoking status (current or former). We then
379 weight each selected participant by his/her inverse probability of selection to fully represent
380 the cohorts of participants who ever smoked at the time of enrollment. To maximize
381 statistical power, we included the 4 cohorts from the targeted discovery phase again in the
382 validation phase, analyzing the same cases as in the targeted discovery phase but selecting
383 1 new subcohort representative per case. Then, for the 10 cohorts that are included for the
384 first time in the validation phase, we selected 2 subcohort representatives per case.

385

386 The validation phase samples will be assayed for absolute quantification of the 21 proteins
387 on the INTEGRAL panel. The cohorts will be divided into training and testing sets (**Table 1**).

388 To maintain full independence of the testing set, the 4 cohorts that contributed to the
389 targeted discovery phase will be included in the training set. In addition to these 4 cohorts,
390 the training set will additionally include the Alpha-Tocopherol, Beta-Carotene Cancer
391 Prevention Study, the Campaign Against Cancer and Heart Disease, the Physicians' Health
392 Study, and the first blood draw from the Women's Health Initiative. The testing set will
393 include the Golestan Cohort Study, the New York University Women's Health Study, the
394 Shanghai Cohort Study, the Southern Community Cohort Study, the Shanghai Men's Health
395 Study, the second blood draw from the Women's Health Initiative, and the Women's Health
396 Study. These groupings were chosen to balance the training and testing sets by
397 geographical location, US racial/ethnic groups, people who currently or formerly smoked,
398 and lung cancer histological types. For the Women's Health Initiative, two independent
399 groups of participants were selected from two blood draws, and we chose to separate these
400 to achieve a similar balance of current and former smoking cases between the training and
401 testing sets.

402

403 Statistical analyses in the validation phase will use the training set to establish flexible
404 parametric survival models that predict absolute risk of lung cancer over 3 years.³⁰
405 Predictors will include a subset of the 21 proteins from the INTEGRAL panel in addition to
406 demographic, health history, and smoking information. The final model will be evaluated in
407 the testing set to measure its calibration (ratio of observed to expected cases) and
408 discrimination. Discrimination analyses will calculate the area under the receiver-operating
409 curve (AUC) and the sensitivity and specificity of the biomarker model at different thresholds.
410 We will also compare its performance directly to existing definitions of screening eligibility
411 including USPSTF criteria and the PLCOm2012 risk model,¹⁴ where our large sample size
412 will ensure we can detect any AUC differences of clinically meaningful magnitude. A
413 sensitivity analysis will exclude late-stage cases with blood draw close to diagnosis from the
414 dataset.

415

416 Nodule Malignancy project

417 The goal of the Nodule Malignancy project is to identify biomarkers that can differentiate
418 benign versus malignant pulmonary nodules, and the study design is based on the following
419 considerations. First, to focus on the actionable time window while maximizing sample size,
420 we included cases diagnosed up to 5 years following blood draw. For lung cancers
421 diagnosed at the baseline screen, the sample collected at baseline was included. This differs
422 from post-diagnostic samples because all individuals participating in LDCT screening are

423 without cancer diagnosis and mostly asymptomatic at baseline. Second, to maximize
424 statistical power and ensure robust discovery results, we included 4 of the LDCT screening
425 studies in the expanded targeted discovery phase (**Figure 1**). Third, the main comparison
426 group is comprised of individuals with benign nodules who did not develop lung cancer,
427 frequency matched on age at enrollment, age at the abnormal finding, age at blood
428 collection, sex, and follow-up time. When multiple study participants with nodules were
429 available as the matched benign nodule-control, we chose participants with higher estimated
430 probability of nodule malignancy based on the Brock model to increase power for nodules
431 with higher malignancy potential.¹⁹ To examine levels of proteins among nodule-free
432 individuals in the screening-eligible population, we also included one control with no nodule
433 findings per case, frequency matched on age at enrollment, age of blood collection, sex, and
434 follow-up time.

435

436 *Targeted discovery phase*

437 The Nodule Malignancy project used a broad targeted discovery phase. We measured all
438 available panels except the Cell Regulation panel, which did not show any robust
439 associations with lung cancer in the Risk Biomarker project full discovery phase (**Table 2**).
440 We included samples from the Pan-Canadian Early Detection of Lung Cancer Study
441 (PanCan), UK Lung Cancer Pilot Screening Trial (UKLS), International Early Lung Cancer
442 Action Program (IELCAP)-Toronto, and Pamplona-IELCAP (**Table 1**; further details in
443 **Supplementary Table 1**). All samples within each LDCT study were randomly plated
444 regardless of their cancer or nodule status to avoid batch effects by case status.

445

446 Statistical analyses applied multivariable logistic regression for each protein, adjusting for
447 the Brock nodule malignancy score which includes age, sex, family history of lung cancer,
448 emphysema, and nodule size, type, location, count, and spiculation (when available).¹⁹ To
449 select protein markers for the INTEGRAL panel, we are using elastic net penalized
450 regression³¹ and a random-forest-based feature selection approach³² to identify the
451 combination of markers that best predicts nodule malignancy. We will also conduct analyses
452 stratified by time to diagnosis. We will prioritize markers based on selection by either elastic
453 net or random forest and consistency of results across studies.

454

455 *Validation phase*

456 To evaluate the results obtained from the targeted discovery based on relative abundance,
457 we will measure the INTEGRAL panel with absolute quantification in the same set of

458 samples (PanCan, UKLS, IELCAP-Toronto, Pamplona-IELCAP), plus 1 independent study,
459 the Pittsburgh Lung Screening Study (PLuSS). The model will be trained on the 4 original
460 studies, and then evaluated in the PLuSS study. This enables evaluation of the data using
461 absolute quantification of the protein markers (using the same set of studies), as well as
462 external validation of the predictive accuracy (using the independent study).

463

464 **Harmonized databases created within the framework of the INTEGRAL** 465 **Risk Biomarker and Nodule Malignancy projects**

466

467 Risk Biomarker Project

468 One challenge for implementing risk-model-based eligibility for lung cancer screening is the
469 unclear generalizability of risk prediction models in diverse worldwide populations.^{13,14,33} We
470 therefore leveraged the infrastructure from the Risk Biomarker project and the Lung Cancer
471 Cohort Consortium to develop a comprehensive study database for lung cancer incidence
472 and mortality. Our vision is that this database will serve as a key resource for future research
473 on lung cancer. For example, additional epidemiologic studies and development and
474 validation of risk prediction tools will likely be needed to support health authorities in making
475 decisions about lung cancer screening implementation over time in different geographical
476 regions, particularly as the tobacco epidemic evolves.

477

478 The cohorts contributing data on all participants to the LC3 harmonized database include
479 most cohorts in the Risk Biomarker project and some additional cohorts. In total, 24 cohorts
480 have contributed data on nearly 3 million participants (**Table 3**, descriptions in **Supplement**).
481 The years of enrollment range from 1985 to 2010 and geographical regions include North
482 America, Europe, Asia, and Australia. More than 69,000 lung cancer cases have been
483 diagnosed during follow-up, including over 7,600 cases among people who never smoked.

484

485 Details on the eligibility criteria, data collection, and outcome ascertainment for each cohort
486 are provided in the **Supplement** and the list of variables in **Table 4**. The variables were
487 chosen to maximize our ability to calculate risk estimates for existing lung cancer prediction
488 models.^{34,35} A summary of methods for harmonization and imputation is provided in the
489 **Supplement**. An initial analysis in the harmonized dataset compared the performance of
490 lung cancer risk models in the United Kingdom.³⁶

491

492 We have defined a priority to facilitate sharing of the LC3 harmonized database. We are
493 currently establishing a legal and technical infrastructure that will allow investigators outside
494 of the LC3 consortium to request permission to remotely access and analyze the data in a
495 secure computing environment. Available data will include the variables listed in **Table 4**, the
496 metabolomics biomarkers measured in the first project of the LC3,³⁷ and eventually the
497 proteomics biomarkers.

498

499 Nodule Malignancy Project

500 For the Nodule Malignancy project, data from 6 LDCT screening studies were harmonized
501 within the framework of ILCCO. In addition to the 5 LDCT screening studies described
502 above, the National Lung Screening Trial (NLST) is also participating in the Nodule
503 Malignancy project for quantitative imaging analysis. The design of each CT screening
504 program including eligibility and recruitment framework is described in the **Supplement**.

505

506 For quality control, data were systematically checked for missing values, outliers,
507 inadmissible values, aberrant distributions, and internal inconsistencies. All procedures were
508 recorded and a central data dictionary was maintained throughout the process. A total of
509 2,088 cases and 42,940 screened individuals from the 6 LDCT screening studies are
510 included in the harmonized database of screening studies (**Supplementary Table 2**). The
511 variables that are compatible across the screening studies are shown in **Table 4**.

512

513 **Perspectives**

514

515 With the advent of LDCT screening, the potential to substantially reduce lung cancer
516 mortality has vastly expanded, and so has the domain of potential research questions. The
517 current work of the INTEGRAL program aims to address two specific ways in which
518 biomarkers might contribute; namely, to improve the selection of individuals for screening,
519 and to better distinguish between malignant and benign nodules on LDCT images. At the
520 completion of our current work, we anticipate that we will have developed a fit-for-purpose
521 biomarker panel that can be applied in both settings. For pre-screening risk assessment, we
522 will deliver an integrated risk prediction model including the biomarkers on the panel and
523 results of a comprehensive independent validation study of its performance. For nodule
524 discrimination, we will establish an integrated nodule probability model including quantitative
525 radiological features and biomarkers.

526

527 If these steps are successful, important work will remain to implement the INTEGRAL panel
528 in clinical practice. While use of biomarkers in lung cancer screening may have advantages,
529 such as more accurate identification of future cases, there are also potential disadvantages
530 such as the need for a blood draw, delay in obtaining biomarker test results, and financial
531 costs. Specific considerations related to biomarker implementation have been outlined.³⁸ We
532 plan to assess whether repeated measurements of the panel could improve our ability to
533 predict lung cancer risk. Implementation studies will be needed to determine the feasibility
534 and acceptability of this approach in practice. The design of future evaluations will require
535 careful consideration, as we consider it infeasible to evaluate the incremental improvement
536 in performance offered by biomarkers in the setting of a randomized trial. Finally, another
537 future goal might be to identify predictors of lung cancer among people who never smoked.

538

539 It is important to note that many other tools exist or are being developed to refine risk
540 estimation for lung cancer, including both biomarkers and risk prediction models. Another
541 important future direction will be to directly compare the performance of these tools or,
542 where feasible and cost-effective, to integrate them. Comparisons should be made in the
543 same set of samples so that discrimination metrics can be directly compared.

544

545 The INTEGRAL biomarker program represents an ambitious initiative to develop a flexible
546 biomarker tool to improve early lung cancer detection via optimized LDCT screening. With a
547 focus on protein biomarkers, the program spans discovery, panel development, model
548 training and validation – all whilst remaining in an observational framework. The forthcoming
549 results from the validation phase of INTEGRAL will provide a definitive benchmark on the
550 potential for circulating protein biomarkers to improve early detection of lung cancer – and
551 most importantly – whether it is justified to introduce them in a screening scenario to inform
552 who should be screened and how to manage nodules.

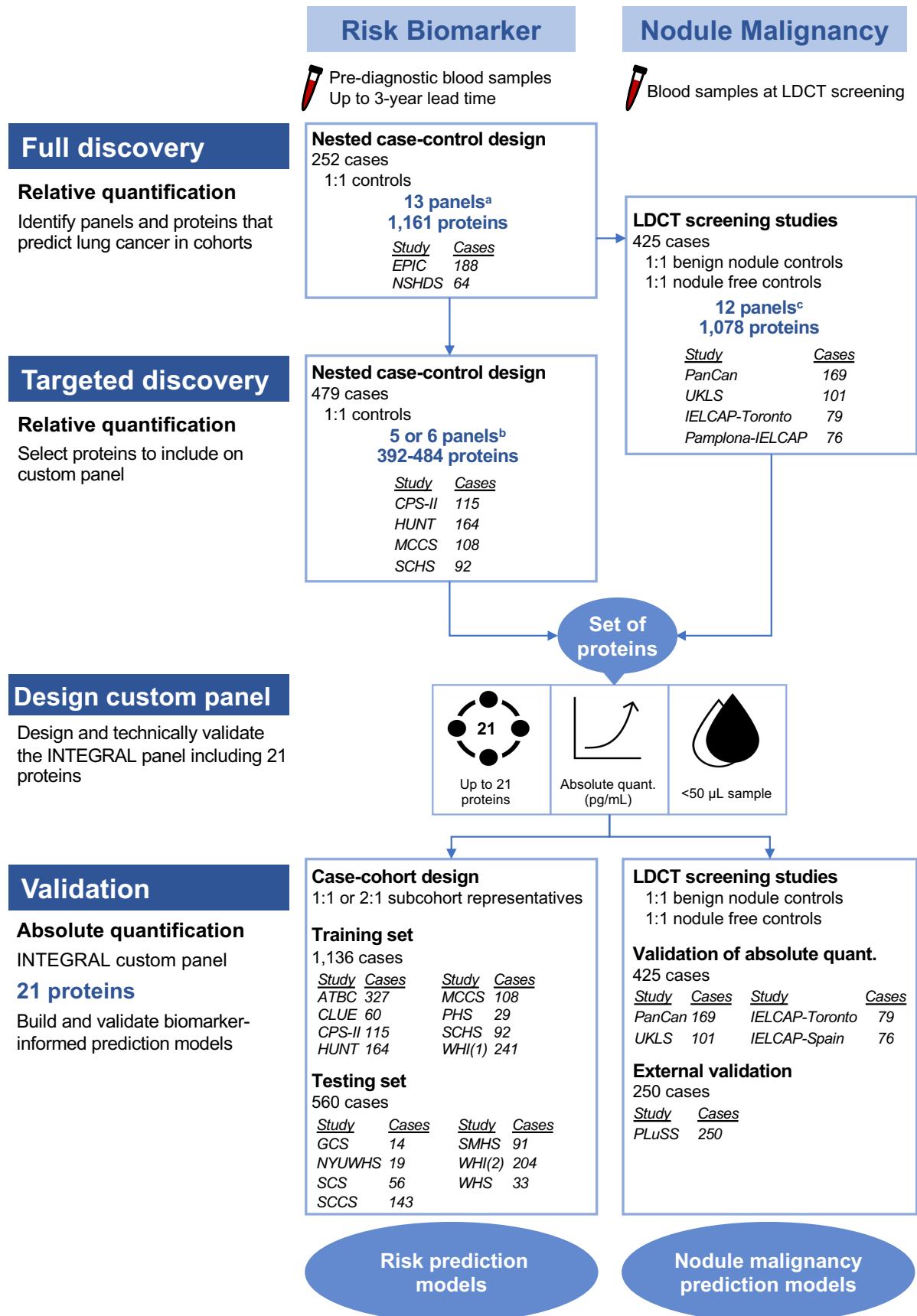
- 555 1. National Lung Screening Trial Research Team, Aberle DR, Adams AM, et al.
556 Reduced lung-cancer mortality with low-dose computed tomographic screening. *N*
557 *Engl J Med*. 2011;365(5):395-409. doi:10.1056/NEJMoa1102873
- 558 2. de Koning HJ, van der Aalst CM, de Jong PA, et al. Reduced lung-cancer mortality
559 with volume CT screening in a randomized trial. *N Engl J Med*. Published online
560 January 29, 2020;10.1056/NEJMoa1911793. doi:10.1056/NEJMoa1911793
- 561 3. Oudkerk M, Liu S, Heuvelmans MA, Walter JE, Field JK. Lung cancer LDCT
562 screening and mortality reduction - evidence, pitfalls and future perspectives. *Nat Rev*
563 *Clin Oncol*. 2021;18(3):135-151. doi:10.1038/s41571-020-00432-6
- 564 4. Fanidi A, Muller DC, Yuan JM, et al. Circulating folate, vitamin B6, and methionine in
565 relation to lung cancer risk in the Lung Cancer Cohort Consortium (LC3). *J Natl*
566 *Cancer Inst*. 2018;110(1). doi:10.1093/jnci/djx119
- 567 5. Muller DC, Hodge AM, Fanidi A, et al. No association between circulating
568 concentrations of vitamin D and risk of lung cancer: an analysis in 20 prospective
569 studies in the Lung Cancer Cohort Consortium (LC3). *Annals of Oncology*.
570 2018;29(6):1468-1475. doi:10.1093/annonc/mdy104
- 571 6. Fanidi A, Carreras-Torres R, Larose TL, et al. Is high vitamin B12 status a cause of
572 lung cancer? *Int J Cancer*. 2019;145(6):1499-1503. doi:10.1002/ijc.32033
- 573 7. Huang JY, Larose TL, Luu HN, et al. Circulating markers of cellular immune activation
574 in prediagnostic blood sample and lung cancer risk in the Lung Cancer Cohort
575 Consortium (LC3). *Int J Cancer*. 2020;146(9):2394-2405. doi:10.1002/ijc.32555
- 576 8. Muller DC, Larose TL, Hodge A, et al. Circulating high sensitivity C reactive protein
577 concentrations and risk of lung cancer: nested case-control study within Lung Cancer
578 Cohort Consortium. *BMJ*. Published online January 3, 2019;k4981.
579 doi:10.1136/bmj.k4981
- 580 9. US National Cancer Institute. NCI Cohort Consortium. Published 2022. Accessed
581 February 21, 2022. <https://epi.grants.cancer.gov/cohort-consortium/>
- 582 10. US Preventive Services Task Force. Screening for Lung Cancer: US Preventive
583 Services Task Force Recommendation Statement. *JAMA*. 2021;325(10):962-970.
584 doi:10.1001/jama.2021.1117
- 585 11. Landy R, Young CD, Skarzynski M, et al. Using prediction models to reduce
586 persistent racial/ethnic disparities in draft 2020 USPSTF lung cancer screening
587 guidelines. *J Natl Cancer Inst*. Published online January 2021.
588 doi:10.1093/jnci/djaa211
- 589 12. Kovalchik SA, Tammemägi M, Berg CD, et al. Targeting of low-dose CT screening
590 according to the risk of lung-cancer death. *N Engl J Med*. 2013;369(3):245-254.
591 doi:10.1056/NEJMoa1301851
- 592 13. Katki HA, Kovalchik SA, Berg CD, Cheung LC, Chaturvedi AK. Development and
593 validation of risk models to select ever-smokers for CT lung cancer screening. *JAMA*.
594 2016;315(21):2300-2311. doi:10.1001/jama.2016.6255
- 595 14. Tammemägi MC, Katki HA, Hocking WG, et al. Selection criteria for lung-cancer
596 screening. *N Engl J Med*. 2013;368(8):728-736. doi:10.1056/NEJMoa1211776
- 597 15. Tammemägi MC, Ruparel M, Tremblay A, et al. USPSTF2013 versus PLCOm2012
598 lung cancer screening eligibility criteria (International Lung Screening Trial): interim
599 analysis of a prospective cohort study. *Lancet Oncol*. Published online December 13,
600 2021. doi:10.1016/S1470-2045(21)00590-8
- 601 16. National Comprehensive Cancer Network. NCCN Clinical Practice Guidelines in
602 Oncology: Lung Cancer Screening version 1.2022. Published 2021. Accessed

- 603 December 20, 2021.
604 https://www.nccn.org/professionals/physician_gls/pdf/lung_screening.pdf
- 605 17. Seijo LM, Peled N, Ajona D, et al. Biomarkers in lung cancer screening:
606 Achievements, promises, and challenges. *Journal of Thoracic Oncology*.
607 2019;14(3):343-357. doi:10.1016/J.JTHO.2018.11.023
 - 608 18. Baldwin D, Callister M, Crosbie PA, et al. Biomarkers in lung cancer screening: the
609 importance of study design. *European Respiratory Journal*. 2021;57(1).
610 doi:10.1183/13993003.04367-2020
 - 611 19. McWilliams A, Tammemagi MC, Mayo JR, et al. Probability of cancer in pulmonary
612 nodules detected on first screening CT. *New England Journal of Medicine*.
613 2013;369(10):910-919. doi:10.1056/NEJMoa1214726
 - 614 20. Al-Ameri A, Malhotra P, Thygesen H, et al. Risk of malignancy in pulmonary nodules:
615 A validation study of four prediction models. *Lung Cancer*. 2015;89(1):27-30.
616 doi:10.1016/j.lungcan.2015.03.018
 - 617 21. Horeweg N, van Rosmalen J, Heuvelmans MA, et al. Lung cancer probability in
618 patients with CT-detected pulmonary nodules: a prespecified analysis of data from the
619 NELSON trial of low-dose CT screening. *Lancet Oncol*. 2014;15(12):1332-1341.
620 doi:10.1016/S1470-2045(14)70389-4
 - 621 22. Feng Z, Pepe MS. Adding rigor to biomarker evaluations—EDRN experience. *Cancer*
622 *Epidemiology Biomarkers & Prevention*. 2020;29(12):2575 LP - 2582.
623 doi:10.1158/1055-9965.EPI-20-0240
 - 624 23. Pepe MS, Feng Z, Janes H, Bossuyt PM, Potter JD. Pivotal evaluation of the
625 accuracy of a biomarker used for classification or prediction: standards for study
626 design. *J Natl Cancer Inst*. 2008;100(20):1432-1438. doi:10.1093/jnci/djn326
 - 627 24. Integrative Analysis of Lung Cancer Etiology and Risk (INTEGRAL) Consortium for
628 Early Detection of Lung Cancer, Guida F, Sun N, et al. Assessment of lung cancer
629 risk on the basis of a biomarker panel of circulating proteins. *JAMA Oncol*. Published
630 online July 12, 2018:e182078. doi:10.1001/jamaoncol.2018.2078
 - 631 25. Larose TL, Meheus F, Brennan P, Johansson M, Robbins HA. Assessment of
632 biomarker testing for lung cancer screening eligibility. *JAMA Netw Open*.
633 2020;3(3):e200409. doi:10.1001/jamanetworkopen.2020.0409
 - 634 26. Silvestri GA, Tanner NT, Kearney P, et al. Assessment of plasma proteomics
635 biomarker's ability to distinguish benign from malignant lung nodules: Results of the
636 PANOPTIC (Pulmonary Nodule Plasma Proteomic Classifier) Trial. *Chest*.
637 2018;154(3):491-500. doi:10.1016/J.CHEST.2018.02.012
 - 638 27. Ostrin EJ, Bantis LE, Wilson DO, et al. Contribution of a blood-based protein
639 biomarker panel to the classification of indeterminate pulmonary nodules. *Journal of*
640 *Thoracic Oncology*. 2021;16(2):228-236. doi:10.1016/j.jtho.2020.09.024
 - 641 28. Olink Proteomics. Measuring protein biomarkers with Olink - technical comparisons
642 and orthogonal validation. Published online 2020. Accessed October 17, 2022.
643 [https://www.olink.com/content/uploads/2021/09/olink-technical-comparisons-and-](https://www.olink.com/content/uploads/2021/09/olink-technical-comparisons-and-orthogonal-validation-1118-v2.0.pdf)
644 [orthogonal-validation-1118-v2.0.pdf](https://www.olink.com/content/uploads/2021/09/olink-technical-comparisons-and-orthogonal-validation-1118-v2.0.pdf)
 - 645 29. Olink Proteomics. Development and validation of customized PEA biomarkers with
646 clinical utility. Published 2017. Accessed October 17, 2022.
647 [https://www.olink.com/content/uploads/2021/09/olink-development-and-validation-of-](https://www.olink.com/content/uploads/2021/09/olink-development-and-validation-of-customized-pea-biomarker-panels-1083-v2.0.pdf)
648 [customized-pea-biomarker-panels-1083-v2.0.pdf](https://www.olink.com/content/uploads/2021/09/olink-development-and-validation-of-customized-pea-biomarker-panels-1083-v2.0.pdf)
 - 649 30. Royston P, Parmar MKB. Flexible parametric proportional-hazards and proportional-
650 odds models for censored survival data, with application to prognostic modelling and
651 estimation of treatment effects. *Stat Med*. 2002;21(15):2175-2197.
652 doi:10.1002/sim.1203

- 653 31. Zou H, Hastie T. Regularization and variable selection via the elastic net. *J R Stat Soc*
654 *Series B Stat Methodol.* 2005;67(2):301-320. doi:[https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-9868.2005.00503.x)
655 9868.2005.00503.x
- 656 32. Degenhardt F, Seifert S, Szymczak S. Evaluation of variable selection methods for
657 random forests and omics data sets. *Brief Bioinform.* 2019;20(2):492-503.
658 doi:10.1093/bib/bbx124
- 659 33. Bach PB, Kattan MW, Thornquist MD, et al. Variations in lung cancer risk among
660 smokers. *J Natl Cancer Inst.* 2003;95(6):470-478.
661 <http://www.ncbi.nlm.nih.gov/pubmed/12644540>
- 662 34. Katki HA, Petito LC, Cheung LC, et al. Implications of 9 risk prediction models for
663 selecting ever-smokers for CT lung-cancer screening. *Ann Intern Med.*
664 2018;169(1):10-19.
- 665 35. Cheung LC, Berg CD, Castle PE, Katki HA, Chaturvedi AK. Life-gained-based versus
666 risk-based selection of smokers for lung cancer screening. *Ann Intern Med.*
667 2019;171(9):623-632. doi:10.7326/M19-1263
- 668 36. Robbins HA, Alcalá K, Swerdlow AJ, et al. Comparative performance of lung cancer
669 risk models to define lung screening eligibility in the United Kingdom. *Br J Cancer.*
670 2021;124(12):2026-2034. doi:10.1038/s41416-021-01278-0
- 671 37. Zahed H, Johansson M, Ueland PM, et al. Epidemiology of 40 blood biomarkers of
672 one-carbon metabolism, vitamin status, inflammation, and renal and endothelial
673 function among cancer-free older adults. *Sci Rep.* 2021;11(1):13805.
674 doi:10.1038/s41598-021-93214-8
- 675 38. Hung RJ. Biomarker-Based Lung Cancer Screening Eligibility: Implementation
676 Considerations. *Cancer Epidemiology, Biomarkers & Prevention.* 2022;31(4):698-701.
677 doi:10.1158/1055-9965.EPI-22-0099
678
679

680
681
682

Figure 1: Schematic describing the development and validation of the INTEGRAL protein panel for lung cancer early detection and nodule malignancy



683

684

685 See Table 1 for definitions of the cohort abbreviations.

686 a: Cardiometabolic, Cardiovascular II, Cardiovascular III, Cell Regulation, Development, Immune response,
687 Inflammation, Metabolism, Neurology Oncology II, Oncology III, Organ Damage, NeuroExploratory

688 b: Cardiovascular III, Inflammation, Immuno-Oncology, Oncology II, Oncology III, NeuroExploratory

689 c: Cardiometabolic, Cardiovascular II, Cardiovascular III, Development, Immune Response, Inflammation,
690 Metabolism, Neurology Oncology II, Oncology III, Organ Damage, NeuroExploratory

Table 1: Description of lung cancer cases participating in the development and validation of the INTEGRAL protein panel for lung cancer early detection and nodule malignancy

Study component	Location	Years of blood draw(s)	Lung cancer cases			Matched controls	Subcohort reps.
			Total	Former smoking	Current smoking		
Risk Biomarker: Full discovery							
European Prospective Investigation into Cancer and Nutrition (EPIC)	Europe	1991-2002	188	59 (31%)	129 (69%)	188	--
Northern Sweden Health and Disease Study (NSHDS)	Sweden	1988-2016	64	26 (41%)	38 (59%)	64	--
Total			252	85 (34%)	167 (66%)	252	
Risk Biomarker: Targeted discovery*							
Cancer Prevention Study II (CPS-II)	USA	1998-2001	115	94 (82%)	21 (18%)	115	--
Nord-Trøndelag Health Study (HUNT)	Norway	1995-1997 2006-2008	164	61 (37%)	103 (63%)	164	--
Melbourne Collaborative Cohort Study (MCCS)**	Australia	1990-1994 2003-2007	108	65 (60%)	43 (40%)	108	--
Singapore Chinese Health Study (SCHS)	Singapore	1994-2005	92	29 (32%)	63 (68%)	92	--
Total			479	249 (52%)	230 (48%)	479	
Risk Biomarker: Validation – training set*							
Alpha-Tocopherol, Beta-Carotene Cancer Prevention Study (ATBC)	Finland	1985-1988	327	--	327 (100%)	--	654
Campaign Against Cancer and Heart Disease (CLUE)	USA	1989-1989	60	33 (55%)	27 (45%)	--	123
Cancer Prevention Study II (CPS-II)	USA	1998-2001	115	94 (82%)	21 (18%)	--	115
Nord-Trøndelag Health Study (HUNT)	Norway	1995-1997 2006-2008	164	61 (37%)	103 (63%)	--	165
Melbourne Collaborative Cohort Study (MCCS)**	Australia	1990-1994 2003-2007	108	65 (60%)	43 (40%)	--	111
Physicians' Health Study (PHS)	USA	1995-2002	29	20 (69%)	9 (31%)	--	58
Singapore Chinese Health Study (SCHS)	Singapore	1994-2005	92	29 (32%)	63 (68%)	--	92
Women's Health Initiative (WHI) (1)**	USA	1993-1997	241	167 (69%)	74 (31%)	--	482
Total			1136	469 (41%)	667 (59%)		1800
Risk Biomarker: Validation – testing set							
Golestan Cohort Study (GCS)	Iran	2004-2008	14	--	14 (100%)	--	28

New York University Women's Health Study (NYUWHS)	USA	1985-1991	19	7 (37%)	12 (63%)	--	38
Shanghai Cohort Study (SCS)	China	1986-1989	56	8 (14%)	48 (86%)	--	112
Southern Community Cohort Study (SCCS)	USA	2002-2009	143	31 (22%)	112 (78%)	--	292
Shanghai Men's Health Study (SMHS)	China	2001-2006	91	19 (21%)	72 (79%)	--	182
Women's Health Initiative (WHI) (2)**	USA	1998-2002	204	145 (71%)	59 (29%)	--	408
Women's Health Study (WHS)	USA	1993-1996	33	19 (58%)	14 (42%)	--	66
Total			560	229 (41%)	331 (59%)		1126
Study component	Location	Years of blood draw(s)	Lung cancer cases			Nodule-free controls	Benign nodule controls
			Total	Former smoking	Current smoking		
<u>Nodule Malignancy: Targeted discovery</u>							
Pan-Canadian Early Detection of Lung Cancer Study (PanCan)	Canada	2008-2014	169	60 (36%)	109 (64%)	169	169
The UK Lung Cancer Pilot Screening Trial (UKLS)	England	2011-2013	101	41 (41%)	60 (59%)	64	92
The International Early Lung Cancer Action Program (IELCAP-Toronto)	Canada	2003-2019	79	30 (38%)	49 (62%)	89	87
The International Early Lung Cancer Action Program (Pamplona-IELCAP)	Spain	2001-2020	76	29 (38%)	47 (62%)	76	82
Total			425	160 (38%)	265 (62%)	398	430
<u>Nodule Malignancy: Validation</u>							
The Pittsburgh Lung Screening Study (PLuSS)	USA	2002-2016	250	77 (31%)	173 (69%)	250	250

INTEGRAL, the Integrative Analysis of Lung Cancer Etiology and Risk program. IELCAP, the International Early Lung Cancer Action Program. Details on the eligibility criteria, data collection, and outcome ascertainment for each cohort are described in the **Supplement**. Further description of the lung cancer cases is given in **Supplementary Table 1**.

*Cohorts in the Risk Biomarker targeted discovery phase are also included in the validation phase training set and are listed twice in the table.

**For the Risk Biomarker project, in MCCS and WHI, participants were sampled separately at two different blood draws. We chose to include the first WHI blood draw in the training set, and the second blood draw in the testing set, to achieve a similar balance of current and former smoking cases between the two sets. For the stratified selection of subcohort representatives, WHI included a stratification by study arm (observational study or the non-intervention arm of the clinical trial).

Table 2: Proteomics panels tested in the full and targeted discovery phases to develop the INTEGRAL protein panel for lung cancer early detection and nodule malignancy

Cohorts	Risk Biomarker Project						Nodule Malignancy Project			
	Full Discovery		Targeted Discovery				Targeted Discovery			
	EPIC	NSHDS	SCHS	CPS-II	HUNT	MCCS	PanCan	UKLS	IELCAP-Toronto	Pamplona-IELCAP
Number of lung cancer cases	188	64	92	115	163	108	169	101	79	76
Number of panels measured	13	13	5	6	5	6	12	12	12	12
Number of measurements*	1196	1196	460	552	460	552	1104	1104	1104	1104
Number of unique proteins*	1161	1161	394	484	392	484	1078	1078	1078	1078
Proteomics panels										
Cardiovascular III	X	X	X	X	X	X	X	X	X	X
Inflammation	X	X	X	X	X	X	X	X	X	X
Immuno-Oncology	(X)	(X)	X	X	X	X	(X)	(X)	(X)	(X)
Oncology II	X	X	X	X	X	X	X	X	X	X
Oncology III	X	X	X	X		X	X	X	X	X
NeuroExploratory	X	X		X	X	X	X	X	X	X
Cardiometabolic	X	X					X	X	X	X
Cardiovascular II	X	X					X	X	X	X
Cell Regulation	X	X								
Development	X	X					X	X	X	X
Immune Response	X	X					X	X	X	X
Metabolism	X	X					X	X	X	X
Neurology	X	X					X	X	X	X
Organ Damage	X	X					X	X	X	X

*Some proteins are measured on multiple panels. In these cases, we chose a single measurement of each protein for analysis by choosing the one that was measured on more cohorts, and then if needed, the one with the highest variance.

(X): all the proteins from the Immuno-Oncology panel are included on other panels assayed as indicated.

Details of the proteins measured on each panel are provided in the Appendix Table.

Table 3: Description of the harmonized Lung Cancer Cohort Consortium database

Cohort	Location	Years of enrollment	Participants, N	Median follow-up (years)*	Female participants, %	Age at enrollment, median (min-max)	----- Lung cancer cases, N (%) -----			
							Total**	Never smoking	Former smoking	Current smoking
AARP	USA	1995-1996	565,645	15.5	40%	62 (50-71)	28,652	2,124 (8)	15,272 (55)	10,189 (37)
ATBC	Finland	1985-1988	29,133	17.7	0%	57 (49-70)	3,959	-	-	3,959 (100)
CLUE	USA	1989	30,461	29.1	57%	48 (18-101)	762	69 (9)	271 (36)	422 (55)
CPS-II	USA	1992-1993	144,670	13.8	55%	70 (47-90)	3,745	446 (12)	2,519 (67)	778 (21)
CSDLH	Canada	1992-1998	11,189	12.3	49%	62 (23-100)	367	65 (18)	203 (56)	93 (26)
EPIC	Europe	1992-2000	518,112	14.9	71%	51 (19-98)	5,233	610 (12)	1,468 (28)	3,155 (60)
GCS	Iran	2004-2008	50,032	13.0	58%	52 (36-78)	118	53 (45)	4 (3)	61 (52)
GS	UK	2003-2009	106,761	9.6	100%	47 (18-102)	217	57 (29)	87 (44)	52 (27)
HPFS	USA	1986	50,444	25.2	0%	55 (32-81)	1,295	164 (13)	635 (51)	444 (36)
HUNT	Norway	1995-1997	78,941	16.9	53%	48 (19-101)	719	34(5)	167 (24)	504 (71)
MCCS	Australia	1990-1994	41,473	23.1	59%	55 (28-76)	855	139 (16)	377 (44)	338 (40)
NHS	USA	1976	120,617	39.9	100%	43 (29-56)	3,986	383 (10)	489 (12)	3,103 (78)
NYUWHS	USA	1985-1991	14,266	30.0	100%	50 (31-70)	484	77 (18)	166 (38)	194 (44)
PHS	USA	1982	26,338	11.7	0%	65 (50-99)	228	49 (21)	127 (56)	52 (23)
PLCO	USA	1993-2001	154,884	11.9	50%	63 (49-78)	3,827	311 (8)	1,821 (50)	1,551 (42)
SCCS	USA	2002-2009	84,429	11.2	60%	52 (40-79)	1,846	109 (6)	369 (21)	1,316 (73)
SCHS	Singapore	1999-2003	50,962	13.5	57%	63 (46-86)	1,300	393 (30)	267 (21)	640 (49)
SCS	China	1986-1989	18,069	25.3	0%	56 (31-79)	1,098	167 (15)	69 (6)	862 (79)
SMHS	China	2002-2006	61,469	12.2	0%	55 (40-75)	1,164	173 (15)	178 (15)	813 (70)
SWHS	China	1996-2000	79,940	18.1	100%	50 (40-70)	975	898 (92)	12 (1)	65 (7)
UKBB	UK	2006-2010	502,105	12.1	54%	57 (37-73)	4,094	728 (18)	1,764 (44)	1,550 (38)
VITAL	USA	2000-2002	77,118	10.0	52%	62 (50-77)	1,374	110 (8)	782 (58)	450 (34)
WHI	USA	1993-1998	118,749	18.2	100%	64 (49-83)	2,389	415 (18)	1,371 (58)	574 (24)
WHS	USA	1992-1995	39,852	24.1	100%	55 (39-90)	588	91 (15)	200 (34)	297 (51)
Total			2,970,659				69,275	7,665 (11)	28,618 (42)	31,462 (47)

*Follow-up time for lung cancer incidence. Mortality follow-up time may differ.

**Cases with missing smoking status are included in the total, but not the stratified counts, so in some cases the stratified counts may not sum to the total.

Details on the eligibility criteria, data collection, and outcome ascertainment for each cohort are described in the **Supplement**. Time varying variables such as age were assessed as of the time of blood draw, or if blood was not collected, as of enrollment. Participants with a history of lung cancer prior to enrollment were excluded. For CSLDH,

the dataset provided is a case-cohort sample (see **Supplement**). For SCHS, the initial enrollment took place during 1993-1998, but the 1999-2003 follow-up visit was used as the baseline for the LC3 dataset (further information in **Supplement**). For WHI, the data include the observational study and the control arms of the Clinical Trials.

AARP: NIH-AARP Diet and Health Study; ATBC: Alpha-Tocopherol, Beta-Carotene Cancer Prevention Study; CLUE: Campaign Against Cancer and Heart Disease II; CPS-II: American Cancer Society Cancer Prevention Study-II Nutrition Cohort; CSDLH: Canadian Study of Diet, Lifestyle and Health; EPIC: European Prospective Investigation into Cancer and Nutrition; GCS: Golestan Cohort Study; GS: Generations Study; HPFS: Health Professionals Follow-up Study; HUNT2 & HUNT3: Trøndelag Health Study; MCCS: Melbourne Collaborative Cohort Study; NHS: Nurses' Health Study I and II; NYUWHS: New York University Women's Health Study; PHS: Physician's Health Study; PLCO: Prostate, Lung, Colorectal, and Ovarian Cancer Screening Trial; SCCS: Southern Community Cohort Study; SCHS: Singapore Chinese Health Study; SCS: Shanghai Cohort Study; SMHS: Shanghai Men's Health Study; UKBB: UK Biobank; VITAL: VITamins And Lifestyle Study; WHI: Women's Health Initiative; WHS: Women's Health Study.

Table 4: Variables included in the harmonized databases for the Lung Cancer Cohort Consortium (Risk Biomarker project) and LDCT screening studies (Nodule Malignancy project)

Variables included in the harmonized Lung Cancer Cohort Consortium database (Risk Biomarker project)				
Demographic information	Follow-up and outcomes	Smoking	Exposures other than smoking	Personal health history
<ul style="list-style-type: none"> • Age • Sex • Education • Race/ethnicity • Year of enrollment or blood draw • State or region of residence (for USA cohorts) 	<ul style="list-style-type: none"> • Follow-up time for lung cancer and death • Lung cancer diagnosis with TNM stage and histology • Vital status and cause of death, including lung cancer death 	<ul style="list-style-type: none"> • Smoking status • Years smoked • Age at smoking initiation • Age at smoking cessation • Years since cessation • Pack-years smoked • Smoking intensity (cigarettes per day) • Type of tobacco product • Time to first cigarette 	<ul style="list-style-type: none"> • Secondhand smoke exposure • Asbestos exposure • Indoor air pollution (e.g. cookstoves) 	<ul style="list-style-type: none"> • Body mass index • Family history of lung cancer • Personal history of cancer • COPD or emphysema • Asthma • Tuberculosis • Daily cough • Liver or kidney condition • Diabetes • Chronic bronchitis • Hypertension • Stroke • Heart attack or heart disease
Variables included in the harmonized LDCT screening study database (Nodule Malignancy project)				
Demographic information	Follow-up and outcomes	Smoking	Nodule characteristics	Personal health history
<ul style="list-style-type: none"> • Age • Sex • Education • Race/ethnicity • Country 	<ul style="list-style-type: none"> • Follow-up time for lung cancer and death • Lung cancer diagnosis with TNM stage and histology • Vital status and cause of death, including lung cancer death 	<ul style="list-style-type: none"> • Smoking status • Duration of smoking • Age at smoking initiation • Age at smoking cessation • Years since quitting • Pack-years smoked • Smoking intensity (cigarettes per day) 	<ul style="list-style-type: none"> • Screening round • Date of screening • Nodule location • Nodule size • Attenuation • Nodule count • Semantic features (spiculation, margin, calcification) • Malignant status 	<ul style="list-style-type: none"> • Body mass index • Family history of lung cancer • Personal history of cancer • COPD • Spirometry measures • Asthma • Chronic bronchitis

Many variables are not available in all cohorts. Cohorts participating in the Risk Biomarker project (see **Table 1**) also provided information on biospecimens including the year of blood draw, storage temperature, number of freeze-thaw cycles, preprocessing time, and details regarding case/control status or subcohort membership.

