

1 **A Novel Nasal Brush-based Classifier of Asthma Identified by Machine Learning Analysis**
2 **of Nasal RNA Sequence Data**

3

4 **Authors:** Gaurav Pandey¹, Om P. Pandey¹, Angela J. Rogers², Gabriel E. Hoffman¹, Benjamin
5 A. Raby³, Scott T. Weiss³, Eric E. Schadt¹, Supinda Bunyavanich^{1,4*}

6

7 **Affiliations:**

8 ¹ Icahn Institute for Genomics and Multiscale Biology and Department of Genetics and Genomic
9 Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

10 ² Division of Pulmonary and Critical Care Medicine, Department of Medicine, Stanford University
11 School of Medicine, Stanford, CA, USA.

12 ³ Channing Division of Network Medicine and Division of Pulmonary and Critical Care Medicine,
13 Brigham & Women's Hospital, and Harvard Medical School, Boston, MA, USA.

14 ⁴ Division of Allergy & Immunology, Department of Pediatrics, Icahn School of Medicine at Mount
15 Sinai, New York, NY, USA.

16

17 *To whom correspondence should be addressed: Supinda Bunyavanich, MD, MPH, Icahn
18 School of Medicine at Mount Sinai, 1425 Madison Avenue #1498, New York, NY 10029, USA,
19 Tel. +1 212 659 8262, Fax +1 212 426 1902, supinda@post.harvard.edu

20 **ABSTRACT:**

21 Asthma is a common, under-diagnosed disease affecting all ages. We sought to identify a nasal
22 brush-based classifier of mild/moderate asthma. One hundred ninety subjects with
23 mild/moderate asthma and controls underwent nasal brushing and RNA sequencing of nasal
24 samples. A machine learning-based pipeline, comprised of feature selection, classification, and
25 statistical analyses, identified a diagnostic classifier of asthma consisting of 90 nasally
26 expressed genes interpreted via an L2-regularized logistic regression classification model. This
27 nasal brush-based classifier performed with strong predictive value and sensitivity across eight
28 validation test sets, including (1) a test set of independent asthmatic and non-asthmatic subjects
29 profiled by RNA sequencing (positive and negative predictive values of 1.00 and 0.96,
30 respectively; AUC of 0.994), (2) two independent case-control cohorts of asthma profiled by
31 microarray, and (3) five independent cohorts of subjects with other respiratory conditions
32 (allergic rhinitis, upper respiratory infection, cystic fibrosis, smoking), where the panel had a low
33 to zero rate of misclassification. Translational development of this classifier into a diagnostic
34 nasal brush-based biomarker for clinical use could aid in asthma detection and care.

35 **Introduction**

36 Asthma is a chronic respiratory disease that affects 8.6% of children and 7.4% of adults
37 in the United States [1]. Its true prevalence may be higher. The fluctuating airflow obstruction,
38 bronchial hyper-responsiveness, and airway inflammation that characterize mild to moderate
39 asthma can be difficult to detect in busy, routine clinical settings [2]. In one study of US middle
40 school children, 11% reported physician-diagnosed asthma with current symptoms, while an
41 additional 17% reported active asthma-like symptoms without a diagnosis of asthma [3].
42 Undiagnosed asthma leads to missed school and work, restricted activity, emergency
43 department visits, and hospitalizations [3, 4]. Given the high prevalence of asthma and
44 consequences of missed diagnosis, there is high potential impact of improved diagnostic tools
45 for asthma [5].

46 National and international guidelines recommend that the diagnosis of asthma should be
47 based on a history of typical symptoms *and objective findings* of variable expiratory airflow
48 limitation [6, 7]. However, obtaining such objective findings can be challenging given currently
49 available tools. Pulmonary function tests (PFTs) require equipment, expertise, and experience
50 to execute well [8, 9]. Many individuals have difficulty with PFTs because they require
51 coordinated breaths into a device. Results are unreliable if the procedure is done with poor
52 technique [8]. Further, PFTs are usually not immediately available in primary care settings.
53 Despite guidelines recommending objective tests such as PFTs to assess possible asthma,
54 PFTs are not done in over half of patients suspected of having asthma [8]. Induced sputum and
55 exhaled nitric oxide have been explored as asthma biomarkers, but their implementation
56 requires technical expertise and does not yield better clinical results than physician-guided
57 management alone [10]. Given the above, the reality is that most asthma is still clinically
58 diagnosed and managed based on self-report [8, 9]. This is suboptimal for mild/moderate
59 asthma given its waxing/waning nature, and because self-reported symptoms and medication
60 use are biased [11].

61 A nasal biomarker of asthma is of high interest given the accessibility of the nose and
62 shared airway biology between the upper and lower respiratory tracts [12-15]. The easily
63 accessible nasal passages are directly connected to the lungs and exposed to common
64 environmental and microbial factors. In this study, we applied next-generation sequencing and
65 machine learning to identify a novel nasal brush-based classifier of asthma (**Figure 1**).
66 Specifically, we used RNA sequencing (RNAseq) to comprehensively profile gene expression
67 from nasal brushings collected from subjects with mild to moderate asthma and controls,
68 creating the largest nasal RNAseq data set in asthma to date. Using a robust machine learning-
69 based pipeline comprised of feature selection [16], classification [17], and statistical analyses
70 [18], we identified an asthma gene panel that accurately differentiates subjects with and without
71 mild-moderate asthma. This pipeline was designed with a systems biology-based perspective
72 that many genes, even ones with marginal effects, can collectively classify phenotypes (here
73 asthma) more accurately than individual genes [19].

74 We validated this asthma gene panel on eight test sets of independent subjects with
75 asthma and other respiratory conditions, finding that it performed with high accuracy, sensitivity,
76 and specificity. As the study of nasal transcriptomics in asthma has been marked by small
77 studies thus far, our relatively large study importantly adds RNAseq data to the field while also
78 leveraging smaller existing data sets for external validation. We see our identification of a
79 diagnostic nasal brush-based classifier of asthma as the first step in the development of
80 minimally invasive, nasal biomarkers for asthma care, with translational development for clinical
81 implementation to follow next. As with any disease, the first step is to accurately identify affected
82 patients, and a next phase of research will be to develop nasal biomarkers to predict treatment
83 response.

84

85 **Results**

86 *Study population and baseline characteristics*

87 We performed nasal brushing on 190 subjects for this study, including 66 subjects with
88 well-defined mild to moderate persistent asthma (based on symptoms, medication need, and
89 demonstrated airway hyper-responsiveness by methacholine challenge) and 124 subjects
90 without asthma (based on no personal or family history of asthma, normal spirometry, and no
91 bronchodilator response). The definitional criteria we used for mild-moderate asthma are
92 consistent with US National Heart Lung Blood Institute guidelines for the diagnosis of asthma
93 [7], and are the same criteria used in the longest NIH-sponsored study of mild-moderate asthma
94 [20, 21].

95 From these 190 subjects, a random selection of 150 subjects were *a priori* assigned as
96 the development set (to be used for asthma classifier development), and the remaining 40
97 subjects were *a priori* assigned as the RNAseq test set (to be used as one of 8 validation test
98 sets for testing of the asthma classifier identified from the development set).

99 The baseline characteristics of the subjects in the development set (n=150) are shown in
100 the left section of **Table 1**. The mean age of subjects with asthma was somewhat lower than
101 subjects without asthma, with slightly more male subjects with asthma and more female
102 subjects without asthma. Caucasians were more prevalent in subjects without asthma, which
103 was expected based on the inclusion criteria. Consistent with reversible airway obstruction that
104 characterizes asthma [2], subjects with asthma had significantly greater bronchodilator
105 response than control subjects (T-test $P = 1.4 \times 10^{-5}$). Allergic rhinitis was more prevalent in
106 subjects with asthma (Fisher's exact test $P = 0.005$), consistent with known comorbidity
107 between allergic rhinitis and asthma [22]. Rates of smoking between subjects with and without
108 asthma were not significantly different.

109 RNA isolated from nasal brushings from the subjects was of good quality, with mean RIN
110 7.8 (± 1.1). The median number of paired-end reads per sample from RNA sequencing was 36.3
111 million. Following pre-processing (normalization and filtering) of the raw RNASeq data, 11,587
112 genes were used for statistical and machine learning analysis. VariancePartition analysis [23],

113 which is designed to analyze the contribution of technical and biological factors to variation in
114 gene expression, showed that age, race, and sex contributed minimally to total gene expression
115 variance (**Supplementary Figure 1**). For this reason, we did not adjust the pre-processed
116 RNASeq data for these factors.

117 Differential gene expression analysis by DeSeq2 [24] showed that 1613 and 1259 genes
118 were respectively over- and under-expressed in asthma cases versus controls (false discovery
119 rate (FDR) ≤ 0.05) (**Supplementary Table 1**). These genes were enriched for disease-relevant
120 pathways in the Molecular Signature Database [25], including immune system (fold
121 change=3.6, FDR= 1.07×10^{-22}), adaptive immune system (fold change=3.91, FDR= 1.46×10^{-15}),
122 and innate immune system (fold change=4.1, FDR= 4.47×10^{-9}) (**Supplementary Table 1**).

123

124 *Identifying a nasal brush-based classifier to predict asthma status*

125 To identify a nasal brush-based classifier that accurately predicts asthma status using
126 the RNAseq data generated, we developed a rigorous machine learning pipeline that combined
127 feature (gene) selection [16] and classification techniques [17] that was applied to the
128 development set (**Materials and Methods** and **Supplementary Figure 2**). This pipeline was
129 designed with a systems biology-based perspective that many genes, even ones with marginal
130 effects, can collectively classify phenotypes (here asthma) more accurately than individual
131 genes. Each gene expression trait can be evaluated on its own or in combination with other
132 gene expression traits to assess how well it distinguishes asthma cases from controls (a
133 process referred to as feature selection). Once the most predictive gene expression traits
134 (features) are identified, various machine learning algorithms can be applied to build a classifier
135 that is optimized to predict asthma status as accurately as possible given the data (a process
136 referred to as classification analysis).

137 Feature selection in our pipeline involved a cross validation-based protocol [26] using
138 the well-established Recursive Feature Elimination (RFE) algorithm [16] combined with L_2 -

139 regularized Logistic Regression (LR or Logistic) and Support Vector Machine (SVM-Linear
140 (kernel)) algorithms [17] (combinations referred to as LR-RFE and SVM-RFE respectively)
141 **(Supplementary Figure 3)**. Classification analysis was then performed by applying four global
142 classification algorithms (SVM-Linear, AdaBoost, Random Forest, and Logistic) [17] to the
143 expression profiles of the gene sets identified by feature selection. To reduce the potential
144 adverse effect of overfitting, this process (feature selection and classification) was repeated 100
145 times on 100 random splits of the development set into training and holdout sets. The final
146 classifier was selected by statistically comparing the models in terms of both classification
147 performance and parsimony, i.e., the number of genes included in the model [18]
148 **(Supplementary Figure 4)**.

149 Due to the imbalance of the two classes (asthma and controls) in our cohort (consistent
150 with imbalances in the general population), we used F-measure as the main evaluation metric in
151 our study [27]. This class-specific measure is a conservative mean of precision (predictive
152 value) and recall (same as sensitivity), and is described in detail in **Box 1** and **Supplementary**
153 **Figure 5**. F-measure can range from 0 to 1, with higher values indicating superior classification
154 performance. An F-measure value of 0.5 does not represent a random model. To provide
155 context for our performance assessments, we also computed commonly used evaluation
156 measures, including positive and negative predictive values (PPVs and NPVs) and Area Under
157 the Receiver Operating Characteristic (ROC) Curve (AUC) scores (**Box 1** and **Supplementary**
158 **Figure 5**).

159

160 **Box 1: Evaluation measures for predictive models**

161 Many measures exist for evaluating the performance of classifiers. The most commonly
162 used evaluation measures in medicine are the positive and negative predictive values (PPV and
163 NPV respectively; **Supplementary Figure 5**), and Area Under the Receiver Operating
164 Characteristic (ROC) Curve (AUC score) [27]. However, these measures have several
165 limitations. PPV and NPV ignore the critical dimension of sensitivity [27]. For instance, a
166 classifier may predict perfectly for only one asthma sample in a cohort and make no predictions
167 for all other asthma samples. This will yield a PPV of 1, but poor sensitivity, since none of the
168 other asthma samples were identified by the classifier. ROC curves and their AUC scores do
169 not accurately reflect performance when the number of cases and controls in a sample are
170 imbalanced [27], which is frequently the case in clinical studies and medical practice. For such
171 situations, precision, recall, and F-measure (**Supplementary Figure 5**) are considered more
172 meaningful performance measures for classifier evaluation. Note that precision for cases (e.g.
173 asthma) is equivalent to PPV, and precision for controls (e.g. no asthma) is equivalent to NPV
174 (**Supplementary Figure 5**). Recall is the same as sensitivity. F-measure is the harmonic
175 (conservative) mean of precision and recall that is computed separately for each class, and thus
176 provides a more comprehensive and reliable assessment of model performance for cohorts with
177 unbalanced class distributions. Like PPV, NPV and AUC, F-measure ranges from 0 to 1, with
178 higher values indicating superior classification performance, but a value of 0.5 for F-measure
179 does not represent a random model and could in some cases indicate superior performance
180 over random. For the above reasons, we consider F-measure as the primary evaluation
181 measure in our study, although we also provide PPV, NPV and AUC measures for context.

182

183 The best performing and most parsimonious combination of feature selection and
184 classification algorithm identified by our machine learning pipeline was LR-RFE & Logistic
185 (Regression) (**Supplementary Figure 4**). The classifier inferred using this combination was built
186 on 90 predictive genes and will be henceforth referred to as the *asthma gene panel*. We
187 emphasize that the expression values of the panel's 90 genes must be used in combination with
188 the Logistic classifier and the model's optimal classification threshold (i.e. predicted
189 label=asthma if classifier's probability output \geq 0.76, else predicted label=no asthma) to be used
190 effectively for asthma classification.

191

192 *Validation of the asthma gene panel classifier in an RNAseq test set of independent subjects*

193 Our next step was to validate the asthma gene panel in an RNAseq test set of
194 independent subjects, for which we used the test set (n=40) of nasal RNAseq data from
195 independent subjects. The baseline characteristics of the subjects in this test set are shown in
196 the right section of **Table 1**. Subjects in the development and test sets were generally similar,
197 except for a lower prevalence of allergic rhinitis among those without asthma in the test set.

198 The asthma gene panel performed with high accuracy in the RNAseq test set's
199 independent subjects, achieving AUC = 0.994 (**Figure 2**), PPV 1.00, and NPV 0.96 (**Figures 3B**
200 **and 3D, left most bar**). In terms of the F-measure metric, the panel achieved F = 0.98 and 0.96
201 for classifying asthma and no asthma, respectively (**Figures 3A and 3C, left most bar**). For
202 comparison, the much lower performance of permutation-based random models is shown in
203 **Supplementary Figure 6**.

204 Our machine learning pipeline evaluated models from several combinations of feature
205 selection and classification algorithms to select the most predictive classifier. Potentially
206 predictive genes can also be identified from differential expression analysis and results from
207 prior asthma-related studies. **Figure 4** shows the performance of the asthma gene panel in the
208 RNAseq test set relative to that of alternative classifiers trained on the development set using:

209 (1) other classifiers tested in our machine learning pipeline, (2) all genes in our data set (11587
210 genes after filtering), (3) all differentially expressed genes in the development set (2872 genes)
211 (**Supplementary Table 1**), (4) genes associated with asthma from prior studies[28] (70 genes)
212 (**Supplementary Table 2**), and (5) a commonly used one-step classification model (L1-Logistic)
213 [29] (243 genes). The asthma gene panel identified by our pipeline outperformed all these
214 alternative classifiers despite its reliance on a small number of genes.

215 We emphasize that our panel produced more accurate predictions than models using all
216 genes, all differentially expressed genes, and all known asthma genes. This supports that data-
217 driven methods can build more effective classifiers than those built exclusively on traditional
218 statistical methods (which do not necessarily target classification), and current domain
219 knowledge (which may be incomplete and subject to investigation bias). Our panel also
220 outperformed and was more parsimonious than the model learned using the commonly used
221 L1-Logistic method, which combined feature selection and classification into a single step. The
222 fact that our asthma gene panel performed well in an independent RNAseq test set while also
223 outperforming alternative models lends confidence to the panel's classification ability.

224

225 *Validation of the asthma gene panel in external asthma cohorts*

226 To assess the generalizability of our asthma gene panel for asthma classification in
227 other populations and profiling platforms, we applied the panel to microarray-derived nasal gene
228 expression data generated from independent cohorts of asthmatics and controls : Asthma1
229 (GEO GSE19187)[30] and Asthma2 (GEO GSE46171)[31]. **Supplementary Table 3**
230 summarizes the characteristics of these external, independent case-control cohorts. In general,
231 RNAseq-based predictive models are not expected to translate well to microarray-profiled
232 samples [32, 33]. A major reason is that gene mappings do not perfectly correspond between
233 RNAseq and microarray due to disparities between array annotations and RNAseq gene models

234 [33]. Our goal was to assess the performance of our asthma gene panel despite discordances in
235 study designs, sample collections, and gene expression profiling platforms.

236 The asthma gene panel performed relatively well (**Figure 3 middle bars**) and
237 consistently better than permutation-based random models (**Supplementary Figure 6**) in
238 classifying asthma and no asthma in both the Asthma1 and Asthma2 microarray-based test
239 sets. The panel achieved similar F-measures in the two test sets (**Figures 3A and 3C middle**
240 **bars**), although the PPV and NPV measures were more dissimilar for Asthma2 (PPV 0.93, NPV
241 0.31) than for Asthma1 (PPV 0.61, NPV 0.67) (**Figure 3B and 3D middle bars**). Although the
242 panel's performance was better than its random counterparts for both these test sets, the
243 difference in this performance was smaller for Asthma2. This occurred partially because
244 Asthma2 includes many more asthma cases than controls (23 vs. 5), which is counter to the
245 expected distribution in the general population. In such a skewed data set, it is possible for a
246 random model to yield an artificially high F-measure for asthma by predicting every sample as
247 asthmatic. We verified that this occurred with the random models tested on Asthma2.

248 To assess how the asthma gene panel might perform in a larger external test set , we
249 combined samples from Asthma1 and Asthma2 and performed the evaluation on this combined
250 set. We chose this approach because no single large, external dataset of nasal gene expression
251 in asthma exists, and combining cohorts could yield a joint test set with heterogeneity that
252 partially reflects real-life heterogeneity of asthma. As expected, all the performance measures
253 for this combined test set were intermediate to those for Asthma1 and Asthma2 (**Figure 3 right**
254 **most bars**), and they still outperformed random counterparts of the panel (**Supplementary**
255 **Figure 6**). These results supported that our panel also performs well in a larger and more
256 heterogeneous cohort.

257 Overall, despite the discordance of gene expression profiling platforms, study designs,
258 and sample collection methods, our asthma gene panel performed reasonably well in these
259 external test sets, supporting a degree of generalizability of the panel across platforms and

260 cohorts. Such a translatable result is not frequently observed in genomic medicine research,
261 especially those based on gene expression [34, 35].

262

263 *Specificity of the asthma gene panel: validation in external cohorts with non-asthma respiratory*
264 *conditions*

265 To assess the specificity of our panel, we next sought to determine if it would misclassify
266 as asthma other respiratory conditions with symptoms that overlap with asthma. To this end, we
267 evaluated the performance of the asthma gene panel on nasal gene expression data derived
268 from case-control cohorts with allergic rhinitis (GSE43523) [36], upper respiratory infection
269 (GSE46171) [31], cystic fibrosis (GSE40445) [37], and smoking (GSE8987) [12].

270 **Supplementary Table 4** details the characteristics for these external cohorts with non-asthma
271 respiratory conditions. In three of these five non-asthma cohorts (Allergic Rhinitis, Cystic
272 Fibrosis and Smoking), the panel appropriately produced one-sided classifications, i.e., samples
273 were all appropriately classified as “no asthma.” This is shown by the zero F-measure for the
274 positive (asthma) class (**Figure 5A**) and perfect F-measure for the negative (no asthma) class
275 (**Figure 5C**) obtained by the panel in these cohorts. In other words, the precision for the asthma
276 class (PPV) of our panel was exactly and appropriately zero (**Figure 5B**), and NPV was
277 perfectly 1.00 for these cohorts with non-asthma conditions (**Figures 5D**). The URI day 2 and 6
278 cohorts were slight deviations from these trends, where the panel achieved perfect NPVs of
279 1.00 (**Figure 5D**), but marginally lower F-measure for the “no asthma” class (**Figure 5C**) due to
280 slightly lower than perfect sensitivity. This may have been influenced by common inflammatory
281 pathways underlying early viral inflammation and asthma [38]. Nonetheless, consistent with the
282 other non-asthma test sets, the panel’s misclassification of URI as asthma was rare and
283 substantially less than its random counterpart classifiers (**Supplementary Figure 7**).

284 To assess the asthma gene panel’s performance if presented with a large,
285 heterogeneous collection of non-asthma respiratory conditions reflective of real clinical settings,

286 we aggregated the non-asthma cohorts into a “Combined non-asthma” test set and applied the
287 asthma gene panel. The results included an appropriately zero F-measure for asthma and zero
288 PPV, and F-measure 0.97 for no asthma and NPV 1.00 (**Figure 5, right most bars**). Results
289 from the individual and combined non-asthma test sets collectively support that the asthma
290 gene panel would rarely misclassify other respiratory diseases as asthma.

291

292 *Statistical and Pathway Examination of Genes in the Asthma Gene Panel*

293 An interesting question to ask for a disease classification panel is how does its predictive
294 ability relate to the individual differential expression status of the genes constituting the panel?
295 We found that 46 of the 90 genes included in our panel were differentially expressed (FDR
296 ≤ 0.05), with 22 and 24 genes over- and under-expressed in asthma respectively (**Figure 6,**
297 **Supplementary Table 1**). More generally, the genes in our panel had lower differential
298 expression FDR values than other genes (Kolmogorov-Smirnov statistic=0.289, P-
299 value= 2.73×10^{-37}) (**Supplementary Figure 8**).

300 In terms of biological function, pathway enrichment analysis of our panel’s 90 genes,
301 though statistically limited by the small number of genes, yielded enrichment for pathways
302 including defense response (fold change=2.86, FDR=0.006) and response to external stimulus
303 (fold change=2.50, FDR=0.012). Only four (*C3*, *DEFB1*, *CYFIP2* and *GSTT1*) of the 90 genes
304 are known asthma genes and are functionally involved in complement activation, microbicidal
305 activity, T-cell differentiation, and oxidative stress, respectively [28]. These results suggest that
306 our machine learning pipeline was able to extract information beyond individually differentially
307 expressed or previously known asthma genes, allowing for the identification of a parsimonious
308 panel of genes that collectively enabled accurate asthma classification.

309 **Discussion**

310 We identified a panel of genes expressed in nasal brushings that accurately classifies
311 subjects with mild/moderate asthma from controls. This nasal brush-based panel, consisting of
312 the expression profiles of 90 genes interpreted via a logistic regression classification model,
313 performed with high precision (PPV=1.00 and NPV=0.96) and recall for classifying asthma
314 (AUC=0.994). The performance of the asthma gene panel across independent asthma test sets
315 demonstrates the generalizability of the panel across study populations and two major
316 modalities of gene expression profiling (RNAseq and microarray). Additionally, the panel's low
317 to zero rate of misclassification on external cohorts with non-asthma respiratory conditions
318 supported the specificity of this panel.

319 Our nasal brush-based asthma gene panel is based on the common biology of the upper
320 and lower airway, a concept supported by clinical practice and previous findings [12-15].
321 Clinically, we rely on the united airway by screening for lower airway infections (e.g. influenza,
322 methicillin-resistant *Staphylococcus aureus*) with nasal swabs [39]. Sridhar et al. found that
323 gene expression consequences of tobacco smoking in bronchial epithelial cells were reflected in
324 nasal epithelium [12]. Wagener et al. compared gene expression in the nasal and bronchial
325 epithelia from 17 subjects, finding that 99% of the 33,000 genes tested exhibited no differential
326 expression between the nasal and bronchial epithelia in those with airway disease [13]. In a
327 study of 30 children, Guajardo et al. identified gene clusters with differential expression in nasal
328 epithelium between subjects with exacerbated asthma vs. controls [14]. The above studies were
329 done with small sample sizes and microarray technology. More recently, Poole et al. compared
330 RNAseq profiles of nasal brushings from 10 asthmatic and 10 control subjects to publicly
331 available bronchial transcriptional data, finding correlation ($\rho = 0.87$) between nasal and
332 bronchial transcripts, as well as correlation ($\rho=0.77$) between nasal differential expression and
333 previously observed bronchial differential expression in asthmatics [15]. To our knowledge, our

334 study has generated the largest nasal RNAseq data set in asthma to date and is the first to
335 identify a nasal brush-based classifier of asthma.

336 Although based on only 90 genes, our asthma gene panel classified asthma with greater
337 accuracy than models based on all genes, all differentially expressed genes, and known asthma
338 genes (**Figure 4**). Its superior performance supports that our machine learning pipeline
339 successfully selected a parsimonious set of informative genes that (1) captures more actionable
340 knowledge than traditional differential expression and genetic association analyses, and (2) cuts
341 through the potential noise of genes irrelevant to asthma. These results show that data-driven
342 methods can build more effective classifiers than those built exclusively on current domain
343 knowledge. About half the genes in our asthma gene panel were not differentially expressed at
344 $FDR \leq 0.05$, and as such would not have been examined with greater interest had we only
345 performed traditional differential expression analysis, which is the main analytic approach of
346 virtually all studies of gene expression in asthma. [12-15, 40, 41]. Consistent with basic
347 hypotheses underlying systems biology approaches, our study demonstrated that the asthma
348 gene panel captures signal from differential expression as well as genes below traditional
349 significance thresholds that may still have a contributory role to asthma classification. Only four
350 of the 90 genes (complement component 3 (C3), defensin beta-1 (DEFB1), cytoplasmic FMR1
351 interacting protein (CYFIP2) and glutathione S-transferase theta 1 (GSTT1)) were previously
352 identified to be relevant to asthma by genetic association studies [28].

353 Our asthma gene panel has the potential to be developed into a minimally invasive
354 biomarker to aid asthma diagnosis at clinical frontlines, where time and resources often
355 preclude pulmonary function testing (PFT). Nasal brushing can be performed quickly, does not
356 require machinery for collection, and implementation of our classification model yields a
357 straightforward, binary result of asthma or no asthma. According to the Global Initiative for
358 Asthma and US National Heart Lung Blood Institute, the diagnosis of asthma should be based
359 on a history of typical symptoms *and objective findings* of variable expiratory airflow limitation by

360 PFT [6, 7]. Practically, however, objective measures are often not obtained. Patients with
361 mild/moderate asthma are frequently asymptomatic at the time of exam. PFTs are often not
362 done, with one study showing that over half of 465,866 patients over age 7 years with newly
363 diagnosed asthma had no PFTs performed within a 3.5 year window surrounding diagnosis [8].
364 Clinicians defer PFTs due to lack of equipment, time, and/or expertise to perform and interpret
365 results [8, 9]. Diagnosing asthma based on history alone contributes to its under-diagnosis, as
366 patients with asthma under-perceive and under-report their symptoms [11]. Misdiagnosis of
367 asthma also occurs frequently given overlapping symptoms between asthma and other
368 conditions [42]. Even if PFTs are obtained, spirometric abnormalities in mild/moderate
369 asthmatics are not always present. An objective, accurate diagnostic classifier that is easy to
370 obtain and interpret with minimal effort from the provider and patient could improve asthma
371 diagnostic accuracy so that appropriate management can then be pursued.

372 Implementation of the asthma gene panel could involve clinicians brushing a patient's
373 nose, placing the brush in a prepackaged tube, and submitting the sample for gene expression
374 profiling targeted to the panel. Some platforms allow for direct transcriptional profiling of tissue
375 without an RNA isolation step, avoiding inconveniences associated with direct RNA work [43,
376 44] and yielding comparable results to RNAseq [45]. Bioinformatic interpretation of the output
377 via the logistic regression-based classifier and classification threshold check could be
378 automated, resulting in a determination of asthma or no asthma for the clinician to consider.
379 Gene expression-based diagnostic classifiers are being successfully used in other disease
380 areas, with prominent examples including the commercially available MammaPrint [46] and
381 Oncotype DX [47] for diagnosing/predicting breast cancer phenotypes. These examples from
382 the cancer field demonstrate an existing path for moving a diagnostic gene panel such as ours
383 to clinical use.

384 Because it takes seconds for nasal brushing, an asthma gene panel such as ours may
385 be attractive to time-strapped clinicians, particularly primary care providers at the frontlines of

386 asthma diagnosis. Asthma is frequently diagnosed and treated in the primary care setting [48]
387 where access to PFTs is often not immediately available. Although PFTs yield results without
388 specimen handling, these advantages do not seem to overcome its logistical limitations as
389 evidenced by their low rate of real-life implementation [8, 9]. The direct costs of our panel are
390 likely to be slightly higher than PFTs. Targeted profiling of our 90-gene panel currently costs
391 about \$100 per sample, while PFTs cost about \$80 according to the Medicare Physician Fee
392 Schedule [49]. However, gene expression profiling costs are likely to decrease [50], and
393 implementation of the asthma gene panel could result in cost *savings* if it reduces the under-
394 diagnosis and misdiagnosis of asthma [4]. Undiagnosed asthma leads to costly healthcare
395 utilization worldwide [4], including in the United States, where asthma accounts for \$56 billion in
396 medical costs, lost school and work days, and early deaths [51]. Clinical implementation of our
397 asthma gene panel could identify undiagnosed asthma, leading to its appropriate management
398 before high healthcare costs from unrecognized asthma are incurred. Given the panel's
399 demonstrated specificity, use of our asthma gene panel could also reduce asthma misdiagnosis
400 by correctly providing a determination of "no asthma" in non-asthmatic subjects with conditions
401 often confused with asthma. Clinical benefit from gene-expression based classification has
402 already been seen in the breast cancer field, where use of the 70-gene panel test MammaPrint
403 to guide chemotherapy in a clinical trial leads to a lower 5-year rate of survival without
404 metastasis compared to standard management [46].

405 We recognize that our asthma gene panel did not perform quite as well in the
406 microarray-based vs. RNAseq-based asthma test sets, which was to be expected due to
407 differences in study design and technological factors between RNAseq and microarray profiling.
408 First, the baseline characteristics and phenotyping of the subjects differed. Subjects in the
409 RNAseq test set were adults who were classified as mild/moderate asthmatic or healthy using
410 the same strict criteria as the development set, which required subjects with asthma to have an
411 objective measure of obstructive airway disease (i.e. positive methacholine challenge

412 response). In contrast, subjects in the Asthma1 microarray test set were all children (i.e. not
413 adults) with nasal pathology, as entry criteria included dust mite allergic rhinitis specifically [30]
414 (**Supplementary Table 3**). Subjects from the Asthma2 cohort were adults who were classified
415 as having asthma or healthy based on history. As mentioned, the diagnosis of asthma based on
416 history alone without objective lung function testing can be inaccurate [52]. The phenotypic
417 differences between these test sets alone could explain differences in performance of our
418 asthma gene panel in these test sets. Second, the differential performance may be due to the
419 difference in profiling approach. Gene mappings do not perfectly correspond between RNAseq
420 and microarray due to disparities between array annotations and RNAseq gene models [33].
421 Compared to microarrays, RNAseq quantifies more RNA species and captures a wider range of
422 signal [40]. Prior studies have shown that microarray-derived models can reliably predict
423 phenotypes based on samples' RNAseq profiles, but the converse does not often hold [33].
424 Despite the above limitations, our asthma gene panel performed with reasonable accuracy in
425 classifying asthma in these independent microarray-based test sets. These results support a
426 degree of generalizability of our panel to asthma populations that may be phenotyped or profiled
427 differently.

428 An effective clinical classifier should have good positive and negative predictive value
429 [53]. In our case, if an individual has asthma, the ideal classifier would reliably indicate asthma
430 so that an accurate diagnosis is made, and if an individual does not have asthma, the ideal
431 classifier would indicate “no asthma” so that misdiagnosis does not occur. This was indeed the
432 case with our asthma gene panel, which achieved high positive and negative predictive values
433 of 1.00 and 0.96 respectively in the RNAseq test set. We also tested our asthma gene panel on
434 independent tests sets of subjects with allergic rhinitis, upper respiratory infection, cystic
435 fibrosis, and smoking, and showed that the panel had a low to zero rate of misclassifying other
436 respiratory conditions as asthma (**Figure 5**). These results were particularly notable for allergic
437 rhinitis, a predominantly nasal condition. Although our panel is based on nasal gene expression,

438 and asthma and allergic rhinitis frequently co-occur [22], our panel did not misdiagnose allergic
439 rhinitis as asthma. Although these conclusions are based on relatively small validation sets due
440 to the scarcity of nasal gene expression data in the public domain, the strong performance of
441 our panel gives hope that it will be generalizable and specific in other larger cohorts as well.

442 One of the current limitations of using RNAseq is the cost of processing large number of
443 samples and generating large datasets. Although we have generated one of the largest nasal
444 RNAseq data set in asthma to date, a future direction of this study is to recruit additional cohorts
445 for nasal gene expression profiling and extend validation of our findings in a prospective
446 manner, which will aid in the panel's path to clinical translation. This will also be facilitated by
447 the rapidly falling costs of sequencing technologies [50], especially if done in a targeted manner.
448 We recognize that our development set was from a single center and its baseline characteristics
449 do not characterize all populations. For example, the development set consisted of adults, and
450 our control subjects were largely Caucasian. However, variancePartition analysis demonstrated
451 minimal contribution of age, race, and gender to gene expression variance in our data
452 (**Supplementary Figure 1**). We also find it reassuring that the panel performed reasonably well
453 in multiple external data sets spanning children and adults of varied racial distributions, and with
454 asthma and other respiratory conditions defined by heterogeneous criteria. Subjects with
455 asthma in our development cohort were not all symptomatic at the time of sampling. The fact
456 that the performance of our asthma gene panel does not rely on symptomatic asthma is a
457 strength, as many mild/moderate asthmatics are only sporadically symptomatic given the
458 fluctuating nature of the disease.

459 We see our diagnostic nasal brush-based classifier of asthma as the first step in the
460 development of nasal biomarkers for multiple aspects of asthma care. As with any disease, the
461 first step is to accurately identify affected patients. The asthma gene panel described in this
462 study provides an accurate path to this critical diagnostic step. With a correct diagnosis, an
463 array of existing asthma treatment options can be considered [6]. A next phase of research will

464 be to develop a nasal biomarker to predict endotypes and treatment response, so that asthma
465 treatment can be targeted, and even personalized, with greater efficiency and effectiveness
466 [54].

467 In summary, we applied RNA sequencing and machine learning to identify a panel of
468 genes expressed in nasal brushings that accurately classifies subjects with mild/moderate
469 asthma from controls. This panel performed with accuracy across independent and external test
470 sets, indicating reasonable generalizability across study populations and gene expression
471 profiling modality, as well as specificity to asthma. Our asthma gene panel has the potential to
472 be developed into a clinical biomarker to aid in asthma diagnosis, as it could be quickly obtained
473 by simple nasal brush, does not require machinery for collection, and can be easily interpreted.
474 Technical translation of panel implementation in the clinical environment, as well as prospective
475 trials of its clinical effectiveness as a diagnostic asthma biomarker, are needed next. If further
476 developed and applied to clinical practice, this nasal brush-based asthma gene panel could
477 improve asthma detection and care.

478 **Materials and Methods**

479 *Study design and subjects*

480 Subjects with mild/moderate asthma were a subset of participants of the Childhood
481 Asthma Management Program (CAMP), a multicenter North American study of 1041 subjects
482 with mild to moderate persistent asthma [20, 21]. Findings from the CAMP cohort have defined
483 current practice and guidelines for asthma care and research [21]. Asthma was defined by
484 symptoms ≥ 2 times per week, use of an inhaled bronchodilator \geq twice weekly or use of daily
485 medication for asthma, and increased airway responsiveness to methacholine ($PC_{20} \leq 12.5$
486 mg/ml). The subset of subjects included in this study were CAMP participants who presented for
487 a visit between July 2011 and June 2012 at Brigham and Women's Hospital (Boston, MA), one
488 of the eight study centers for CAMP.

489 Subjects with "no asthma" were recruited during the same time period by advertisement
490 at Brigham & Women's Hospital. Selection criteria were no personal history of asthma, no family
491 history of asthma in first-degree relatives, and self-described Caucasian ethnicity. Participation
492 was limited to Caucasian individuals because a concurrent independent study was planned that
493 would compare these same subjects to 968 Caucasian CAMP subjects who participated in the
494 CAMP Genetics Ancillary study [55]. Subjects underwent pre- and post-bronchodilator
495 spirometry according to American Thoracic Society guidelines. Only those meeting selection
496 criteria and with demonstrated normal lung function without bronchodilator response were
497 considered to have "no asthma."

498

499 *Nasal brushing and RNA sequencing*

500 Nasal brushing was performed with a cytology brush. Brushes were immediately placed
501 in RNALater (ThermoFisher Scientific, Waltham, MA) and then stored at 40°C until RNA
502 extraction. RNA extraction was performed with Qiagen RNeasy Mini Kit (Valencia, CA).

503 Samples were assessed for yield and quality using the 2100 Bioanalyzer (Agilent Technologies,
504 Santa Clara, CA) and Qubit fluorometry (Thermo Fisher Scientific, Grand Island, NY).

505 Of the 190 subjects who underwent nasal brushing (66 with mild/moderate asthma, 124
506 with no asthma), a random selection of 150 subjects were *a priori* assigned as the development
507 set (for classification model development), with the 40 remaining subjects earmarked to serve
508 as a test set of independent subjects (for testing the classification model). To minimize potential
509 batch effects, all samples were submitted together for RNA sequencing (RNAseq). Staff at the
510 Mount Sinai genomics core were blinded to the assignment of samples as development or test
511 set. The sequencing library was prepared with the standard TruSeq RNA Sample Prep Kit v2
512 protocol (Illumina). The mRNA libraries were sequenced on the Illumina HiSeq 2500 platform
513 with a per-sample target of 40-50 million 100 bp paired-end reads. The data were put through
514 Mount Sinai's standard mapping pipeline[56] (using Bowtie [57] and TopHat [58], and
515 assembled into gene- and transcription-level summaries using Cufflinks [59]). Mapped data
516 were subjected to quality control with FastQC and RNA-SeQC [60]. Data were pre-processed
517 separately for the development and test sets to avoid leakage of information across the two data
518 sets and maintain fairness of the machine learning procedures as much as possible. Genes with
519 fewer than 100 counts in at least half the samples were dropped to reduce the potentially
520 adverse effects of noise. DESeq2 [24] was used to normalize the data sets using its variance
521 stabilizing transformation method.

522

523 *VariancePartition Analysis of Potential Confounders*

524 Given differences in age, race, and sex distributions between the asthma and “no
525 asthma” classes, we used the variancePartition method [23] to assess the degree to which
526 these variables influenced gene expression and potentially confounded the target phenotype
527 (asthma status). The total variance in gene expression was partitioned into the variance
528 attributable to age, race, and sex using a linear mixed model implemented in variancePartition

529 v1.0.0 [23]. Age (continuous variable) was modeled as a fixed effect while race and sex
530 (categorical variables) were modeled as random effects. The results showed that age, race, and
531 sex accounted for minimal contributions to total gene expression variance (**Supplementary**
532 **Figure 1**). Downstream analyses were therefore performed with gene expression data
533 unadjusted for these variables.

534

535 *Differential gene expression and pathway enrichment analysis*

536 DESeq2 [24] was used to identify differentially expressed genes in the development set.
537 Genes with $FDR \leq 0.05$ were deemed differentially expressed, with fold change < 1 implying
538 under-expression and vice versa. To identify the functions underlying these genes, pathway
539 enrichment analysis was performed using the Gene Set Enrichment Analysis method applied to
540 the Molecular Signature Database (MSigDB) [25].

541

542 *Identification of the Asthma Gene Panel by Machine Learning Analyses of the RNAseq*

543 *Development Set*

544 To identify gene expression-based classifiers that predict asthma status, we applied a
545 rigorous machine learning pipeline implemented in Python using the scikit-learn package [61]
546 that combined feature (gene) selection [16], classification [17], and statistical analyses of
547 classification performance [18] to the development set (**Supplementary Figure 2**). Feature
548 selection and classification were applied to a training set comprised of 120 randomly selected
549 samples from the development set ($n=150$) as described below. For an independent evaluation
550 of the candidate classifiers generated from the training set by this process, they were then
551 evaluated on the remaining 30 samples (holdout set). Finally, to reduce the dependence of the
552 finally chosen classifier on a specific training-holdout split, this process was repeated 100 times
553 on 100 random splits of the development set into training and holdout sets. The details of the
554 overall process as well as the individual components are as follows.

555 *Feature selection:* The purpose of the feature selection component was to identify
556 subsets of the full set of genes in the development set, whose expression profiles could be used
557 to predict the asthma status as accurately as possible. The two main computations constituting
558 this component were (i) the optimal number of features that should be selected, and (ii) the
559 identification of this number of genes from the full gene set. To reduce the likelihood of
560 overfitting when conducting both these computations on the entire training set, we used a 5x5
561 nested (outer and inner) cross-validation (CV) setup [26] for selecting features from the training
562 set (**Supplementary Figure 3**). The inner CV round was used to determine the optimal number
563 of genes to be selected, and the outer CV round was used to select the set of predictive genes
564 based on this number, thus separating the samples on which these decisions are made. The
565 supervised Recursive Feature Elimination (RFE) algorithm [62] was executed on the inner CV
566 training split to determine the optimal number of features. The use of RFE within this setting
567 enabled us to identify groups of features that are collectively, but not necessarily individually,
568 predictive. This reflects our systems biology-based expectation that many genes, even ones
569 with marginal effects, can play a role in classifying diseases/phenotypes (here asthma) in
570 combination with other more strongly predictive genes [19]. Specifically, we used the L2-
571 regularized Logistic Regression (LR or Logistic) [63] and SVM-Linear (kernel) [64] classification
572 algorithms in conjunction with RFE (combinations henceforth referred to as LR-RFE and SVM-
573 RFE respectively). For this, for a given inner CV training split, all the features (genes) were
574 ranked using the absolute values of the weights assigned to them by an inner classification
575 model, trained using the LR or SVM algorithm, over this split. Next, for each of the conjunctions,
576 the set of top-k ranked features, with k starting with 11587 (all filtered genes) and being reduced
577 by 10% in each iteration until k=1, was considered. The discriminative strength of feature sets
578 consisting of the top k features as per this ranking was assessed by evaluating the performance
579 of the LR or SVM classifier based on them over all the inner CV training-test splits. The optimal
580 number of features to be selected was determined as the value of k that produces the best

581 performance. Next, a ranking of features was derived from the outer CV training split using
582 exactly the same procedure as applied to the inner CV training split. The optimal number of
583 features determined above was selected from the top of this ranking to determine the optimal
584 set of predictive features for this outer CV training split. Executing this process over all the five
585 outer CV training splits created from the development set identified five such sets. Finally, the
586 set of features (genes) that was common to all these sets (i.e. in their intersection/overlap),
587 which is expected to yield a more robust feature set than the individual outer CV splits, was
588 selected as the predictive gene set for this training set. One such set was identified for each of
589 LR-RFE and SVM-RFE.

590 *Classification analyses:* Once predictive gene sets had been selected from feature
591 selection, four global classification algorithms (L2-regularized Logistic Regression (LR or
592 Logistic) [63], SVM-Linear [64], AdaBoost [65], and Random Forest (RF) [66]) were used to
593 learn *intermediate classification models* over the training set. These intermediate models were
594 then applied to the corresponding holdout set to generate probabilistic asthma predictions for
595 the samples. An optimal threshold for converting these probabilistic predictions into binary ones
596 (higher than threshold=asthma, lower than threshold=no asthma) was then computed as the
597 threshold that yielded the highest classification performance on the holdout set. This
598 optimization resulted in the *proposed classification models*.

599 *Statistical analyses of classification performance:* After the above components have
600 been run on 100 training-holdout splits of the development set, we obtain 100 proposed
601 classification models for each of eight feature selection-global classification combinations (two
602 feature selection algorithms (LR-RFE and SVM-RFE) and four global classification algorithms
603 Logistic, SVM-Linear, AdaBoost and RF). The next step of our pipeline was to determine the
604 best performing combination. Instead of making this determination just based on the highest
605 evaluation score, as is typically done in ML studies, we utilized this large population of models
606 and their optimized holdout evaluation scores to conduct a statistical comparison to make this

607 determination. Specifically, we applied the Friedman test followed by the Nemenyi test [18, 67]
608 to this population of modules and their evaluation scores. These tests, which account for
609 multiple hypothesis testing, assessed the statistical significance of the relative difference of
610 performance of the combinations in terms of their relative ranks across the 100 splits.

611 *Optimization for parsimony.* For an effective phenotype classifier, it is essential to
612 consider parsimony in model selection (i.e. minimize number of features (i.e. genes)) to
613 enhance its biological and clinical utility and acceptability. To enforce this for our classifier, an
614 adapted performance measure, defined as the absolute performance measure (F-measure)
615 divided by the number of genes in that model, was used for the above statistical comparison,
616 i.e. as input to the Friedman-Nemenyi tests. In terms of this measure, a model that does not
617 obtain the best performance measure among all models, but uses much fewer genes than the
618 others, may be judged to be the best model. The result of the statistical comparison using this
619 adapted measure was visualized as a Critical Difference plot [18] (**Supplementary Figure 4**),
620 and enabled us to identify the best combination of feature selection and classification method as
621 the left-most entry in this plot.

622 *Final model development:* The final step in our pipeline was to determine the
623 representative model out of the 100 learned the above best combination by finding which of
624 these models yielded the highest evaluation measure (F-measure). In case of ties among
625 multiple candidates, the gene set that produced the best average asthma classification F-
626 measure (**Box 1** and **Supplementary Figure 5**) across all four global classification algorithms
627 was chosen as the gene set constituting the representative model for that combination. This
628 analysis yielded the representative gene set, global classification algorithm, and the optimized
629 asthma classification threshold. Finally, our asthma gene panel was built by training the global
630 classification algorithm to the expression profiles of the representative gene set, and using the
631 optimized threshold for classifying samples with and without asthma.

632

633 *Validation of the Asthma Gene Panel in an RNAseq test set of independent subjects*

634 The asthma gene panel identified by our machine learning pipeline was then tested on
635 the RNAseq test set (n=40) to assess its performance in independent subjects. F-measure was
636 used as the primary measure for classification performance, as described in **Box 1** and
637 **Supplementary Figure 5**. AUC, PPV and NPV were additionally calculated for context.

638

639 *Performance Comparison to Alternative Classification Models*

640 For comparison, the same machine learning methodology was used to train and
641 evaluate models from all combinations of feature selection and global classification methods
642 considered in our pipeline. We also applied our machine learning pipeline with replacement of
643 the feature (gene) selection step with these pre-determined gene sets: (1) all filtered RNAseq
644 genes, (2) all differentially expressed genes, and (3) known asthma genes from a recent review
645 of asthma genetics [28]. To maintain consistency with the machine learning pipeline-derived
646 models, these were each used as a predetermined gene set that was run through the same
647 pipeline (**Supplementary Figure 2** with the feature selection component turned off) to identify
648 the best performing global classification algorithm and the optimal asthma classification
649 threshold for this predetermined set of features. The algorithm and threshold were used to train
650 each of these gene sets' representative classification model over the entire development set,
651 and the resulting model for each of these gene sets was then evaluated on the RNAseq test set.
652 Finally, as a baseline representative of alternative sparse classification algorithms, which
653 represent a one-step option for doing feature selection and classification simultaneously, we
654 also trained an L1-regularized logistic regression model (L1-Logistic) [29] on the development
655 set and evaluated it on the RNAseq test set.

656

657 *Performance Comparison to Permutation-based Random Models*

658 To determine the extent to which the performance of all the above classification models
659 could have been due to chance, we compared their performance with that of their random
660 counterpart models (**Supplementary Figure 6, Supplementary Figure 7**). These counterparts
661 were obtained by randomly permuting the labels of the samples in the development set and
662 executing each of the above model training procedures on these randomized data sets in the
663 same way as for the real development set. These random models were then applied to each of
664 the test sets considered in our study, and their performances were also evaluated in terms of
665 the same measures. For each of real models tested in our study, 100 corresponding random
666 models were learned and evaluated as above, and the performance of the real models was
667 compared with the average performance of the corresponding random models.

668

669 *Validation of the asthma gene panel in external independent asthma cohorts*

670 To assess the generalizability of the asthma gene panel to other populations,
671 microarray-profiled data sets of nasal gene expression from two external asthma cohorts--
672 Asthma1 (GSE19187) [30] and Asthma2 (GSE46171) [31] (**Supplementary Table 3**)-- were
673 obtained from NCBI Gene Expression Omnibus (GEO) [68]. The asthma gene panel was then
674 applied and its performance evaluated on these external asthma cohorts..

675

676 *Validation of the asthma gene panel in external cohorts with other respiratory conditions*

677 To assess the panel's ability to distinguish asthma from respiratory conditions that can
678 have overlapping symptoms with asthma, i.e. its specificity to asthma, microarray-profiled data
679 sets of nasal gene expression were also obtained for five external cohorts with allergic rhinitis
680 (GSE43523) [36], upper respiratory infection (GSE46171) [31], cystic fibrosis (GSE40445) [37],
681 and smoking (GSE8987) [12] (**Supplementary Table 4**). The asthma gene panel was then
682 applied and its performance evaluated on these external cohorts with non-asthma respiratory
683 conditions.

684 **Declarations**

685 *Ethics approval and consent to participate:* The institutional review boards of Brigham &
686 Women's Hospital and the Icahn School of Medicine at Mount Sinai approved the study
687 protocols. Written informed consent was obtained from all subjects.

688

689 *Competing interests:* SB, GP, and EES have filed a patent application related to the findings of
690 this manuscript. The remaining authors declare that they have no competing interests.

691

692 *Funding:* This study was supported by the US National Institutes of Health (NIH R01AI118833,
693 K08AI093538, R01GM114434) and the Icahn Institute for Genomics and Multiscale Biology.

694

695 *Author contributions:* SB directed the study. SB, BAR, and EES designed the study. SB and
696 AJR directed the recruitment of subjects and sample collection. BAR and STW provided
697 guidance for access to subjects. EES advised on sequencing strategy. SB curated the clinical
698 data. SB, GP, and OPP designed and performed the statistical and computational analyses. SB
699 and GP wrote the manuscript. SB, GP, OPP, AJR, GEH, BAR, STW, and EES edited the
700 manuscript. All authors contributed significantly to the work presented in this paper.

701

702 *Acknowledgments:* We thank Kathryn Paul, Laura Ting, Anne Plunkett, Nancy Madden, Ann
703 Fuhlbrigge, Kelan Tantisira, Dan Cossette, Aimee Garciano, and Roxanne Kelly for their
704 assistance and support with recruitment, specimen collection, and sample processing. We thank
705 Robert Griffin and Ana Stanescu for critically reviewing the paper.

References

1. **Current Asthma Prevalence Percents by Age, Sex, and Race/Ethnicity, United States, 2012. Asthma Surveillance Data.** *National Health Interview Survey, National Center for Health Statistics, Centers for Disease Control and Prevention* www.cdc.gov/asthma/asthmadata.htm, downloaded 1/30/2017.
2. Fanta CH: **Asthma.** *N Engl J Med* 2009, **360**:1002-1014.
3. Yeatts K, Shy C, Sotir M, Music S, Herget C: **Health consequences for children with undiagnosed asthma-like symptoms.** *Arch Pediatr Adolesc Med* 2003, **157**:540-544.
4. Stempel DA, Spahn JD, Stanford RH, Rosenzweig JR, McLaughlin TP: **The economic impact of children dispensed asthma medications without an asthma diagnosis.** *J Pediatr* 2006, **148**:819-823.
5. Szeffler SJ, Wenzel S, Brown R, Erzurum SC, Fahy JV, Hamilton RG, Hunt JF, Kita H, Liu AH, Panettieri Jr RA, et al: **Asthma outcomes: Biomarkers.** *Journal of Allergy and Clinical Immunology* 2012, **129**:S9-S23.
6. Reddel HK, Bateman ED, Becker A, Boulet LP, Cruz AA, Drazen JM, Haahtela T, Hurd SS, Inoue H, de Jongste JC, et al: **A summary of the new GINA strategy: a roadmap to asthma control.** *Eur Respir J* 2015, **46**:622-639.
7. **Expert Panel Report 3: Guidelines for the Diagnosis and Management of Asthma.** Washington DC: National Heart Lung and Blood Institute and National Asthma Education and Prevention Program; 2007.
8. Gershon AS, Victor JC, Guan J, Aaron SD, To T: **Pulmonary function testing in the diagnosis of asthma: a population study.** *Chest* 2012, **141**:1190-1196.
9. Sokol KC, Sharma G, Lin YL, Goldblum RM: **Choosing wisely: adherence by physicians to recommended use of spirometry in the diagnosis and management of adult asthma.** *Am J Med* 2015, **128**:502-508.

10. Petsky HL, Cates CJ, Lasserson TJ, Li AM, Turner C, Kynaston JA, Chang AB: **A systematic review and meta-analysis: tailoring asthma treatment on eosinophilic markers (exhaled nitric oxide or sputum eosinophils).** *Thorax* 2012, **67**:199-208.
11. van Schayck CP, van Der Heijden FM, van Den Boom G, Tirimanna PR, van Herwaarden CL: **Underdiagnosis of asthma: is the doctor or the patient to blame? The DIMCA project.** *Thorax* 2000, **55**:562-565.
12. Sridhar S, Schembri F, Zeskind J, Shah V, Gustafson AM, Steiling K, Liu G, Dumas YM, Zhang X, Brody JS, et al: **Smoking-induced gene expression changes in the bronchial airway are reflected in nasal and buccal epithelium.** *BMC Genomics* 2008, **9**:259.
13. Wagener AH, Zwinderman AH, Luiten S, Fokkens WJ, Bel EH, Sterk PJ, van Drunen CM: **The impact of allergic rhinitis and asthma on human nasal and bronchial epithelial gene expression.** *PLoS One* 2013, **8**:e80257.
14. Guajardo JR, Schleifer KW, Daines MO, Ruddy RM, Aronow BJ, Wills-Karp M, Hershey GK: **Altered gene expression profiles in nasal respiratory epithelium reflect stable versus acute childhood asthma.** *J Allergy Clin Immunol* 2005, **115**:243-251.
15. Poole A, Urbanek C, Eng C, Schageman J, Jacobson S, O'Connor BP, Galanter JM, Gignoux CR, Roth LA, Kumar R, et al: **Dissecting childhood asthma with nasal transcriptomics distinguishes subphenotypes of disease.** *J Allergy Clin Immunol* 2014, **133**:670-678 e612.
16. Saeys Y, Inza I, Larranaga P: **A review of feature selection techniques in bioinformatics.** *Bioinformatics* 2007, **23**:2507-2517.
17. Witten IH, Frank E, Hall MA: *Data mining : practical machine learning tools and techniques.* 3rd edn. Burlington, MA: Morgan Kaufmann; 2011.
18. Demsar J: **Statistical Comparisons of Classifiers over Multiple Data Sets.** *J Mach Learn Res* 2006, **7**:1-30.

19. Schadt EE, Friend SH, Shaywitz DA: **A network view of disease and compound screening.** *Nat Rev Drug Discov* 2009, **8**:286-295.
20. **The Childhood Asthma Management Program (CAMP): design, rationale, and methods.** Childhood Asthma Management Program Research Group. *Control Clin Trials* 1999, **20**:91-120.
21. Covar RA, Fuhlbrigge AL, Williams P, Kelly HW, the Childhood Asthma Management Program Research G: **The Childhood Asthma Management Program (CAMP): Contributions to the Understanding of Therapy and the Natural History of Childhood Asthma.** *Curr Respir Care Rep* 2012, **1**:243-250.
22. Egan M, Bunyavanich S: **Allergic rhinitis: the "Ghost Diagnosis" in patients with asthma.** *Asthma Research and Practie* 2015, **1**:DOI: 10.1186/s40733-40015-40008-40730.
23. Hoffman GE, Schadt EE: **variancePartition: Quantifying and interpreting drivers of variation in complex gene expression studies.** *bioRxiv* 2016:doi: <http://dx.doi.org/10.1101/040170>.
24. Love MI, Huber W, Anders S: **Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2.** *Genome Biol* 2014, **15**:550.
25. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP: **Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles.** *Proc Natl Acad Sci U S A* 2005, **102**:15545-15550.
26. Whalen S, Pandey OP, Pandey G: **Predicting protein function and other biomedical characteristics with heterogeneous ensembles.** *Methods* 2016, **93**:92-102.
27. Lever J, Krzywinski M, Altman N: **Points of Significance: Classification Evaluation.** *Nature Methods* 2016, **13**:603-604.

28. Mathias RA: **Introduction to genetics and genomics in asthma: genetics of asthma.** *Adv Exp Med Biol* 2014, **795**:125-155.
29. Vidaurre D, Bielza C, Larrañaga P: **A Survey of L1 Regression.** *International Statistical Review* 2013, **81**:361-387.
30. Giovannini-Chami L, Marcet B, Moreilhon C, Chevalier B, Illie MI, Lebrigand K, Robbe-Sermesant K, Bourrier T, Michiels JF, Mari B, et al: **Distinct epithelial gene expression phenotypes in childhood respiratory allergy.** *Eur Respir J* 2012, **39**:1197-1205.
31. McErlean P, Berdnikovs S, Favoreto S, Jr., Shen J, Biyasheva A, Barbeau R, Easley C, Barczak A, Ward T, Schleimer RP, et al: **Asthmatics with exacerbation during acute respiratory illness exhibit unique transcriptional signatures within the nasal mucosa.** *Genome Med* 2014, **6**:1.
32. Zhang W, Yu Y, Hertwig F, Thierry-Mieg J, Zhang W, Thierry-Mieg D, Wang J, Furlanello C, Devanarayan V, Cheng J, et al: **Comparison of RNA-seq and microarray-based models for clinical endpoint prediction.** *Genome Biol* 2015, **16**:133.
33. Su Z, Fang H, Hong H, Shi L, Zhang W, Zhang W, Zhang Y, Dong Z, Lancashire LJ, Bessarabova M, et al: **An investigation of biomarkers derived from legacy microarray data for their utility in the RNA-seq era.** *Genome Biol* 2014, **15**:523.
34. Venet D, Dumont JE, Detours V: **Most Random Gene Expression Signatures Are Significantly Associated with Breast Cancer Outcome.** *PLoS computational biology* 2011, **7**:e1002240.
35. Chibon F: **Cancer gene expression signatures - the rise and fall?** *Eur J Cancer* 2013, **49**:2000-2009.
36. Imoto Y, Tokunaga T, Matsumoto Y, Hamada Y, Ono M, Yamada T, Ito Y, Arinami T, Okano M, Noguchi E, Fujieda S: **Cystatin SN upregulation in patients with seasonal allergic rhinitis.** *PLoS One* 2013, **8**:e67057.

37. Clarke LA, Sousa L, Barreto C, Amaral MD: **Changes in transcriptome of native nasal epithelium expressing F508del-CFTR and intersecting data from comparable studies.** *Respir Res* 2013, **14**:38.
38. Oliver BG, Robinson P, Peters M, Black J: **Viral infections and asthma: an inflammatory interface?** *Eur Respir J* 2014, **44**:1666-1681.
39. Cowling BJ, Chan KH, Fang VJ, Lau LL, So HC, Fung RO, Ma ES, Kwong AS, Chan CW, Tsui WW, et al: **Comparative epidemiology of pandemic and seasonal influenza A in households.** *N Engl J Med* 2010, **362**:2175-2184.
40. Bunyavanich S, Schadt EE: **Systems biology of asthma and allergic diseases: A multiscale approach.** *J Allergy Clin Immunol* 2014.
41. Sordillo J, Raby BA: **Gene expression profiling in asthma.** *Adv Exp Med Biol* 2014, **795**:157-181.
42. Scott S, Currie J, Albert P, Calverley P, Wilding JP: **Risk of misdiagnosis, health-related quality of life, and BMI in patients who are overweight with doctor-diagnosed asthma.** *Chest* 2012, **141**:616-624.
43. Kulkarni MM: **Digital multiplexed gene expression analysis using the NanoString nCounter system.** *Curr Protoc Mol Biol* 2011, **Chapter 25**:Unit25B 10.
44. Veldman-Jones MH, Brant R, Rooney C, Geh C, Emery H, Harbron CG, Wappett M, Sharpe A, Dymond M, Barrett JC, et al: **Evaluating Robustness and Sensitivity of the NanoString Technologies nCounter Platform to Enable Multiplexed Gene Expression Analysis of Clinical Samples.** *Cancer Res* 2015, **75**:2587-2593.
45. Leong HS, Galletta L, Etemadmoghadam D, George J, Australian Ovarian Cancer S, Kobel M, Ramus SJ, Bowtell D: **Efficient molecular subtype classification of high-grade serous ovarian cancer.** *J Pathol* 2015, **236**:272-277.

46. Cardoso F, van't Veer LJ, Bogaerts J, Slaets L, Viale G, Delalogue S, Pierga JY, Brain E, Causeret S, DeLorenzi M, et al: **70-Gene Signature as an Aid to Treatment Decisions in Early-Stage Breast Cancer.** *N Engl J Med* 2016, **375**:717-729.
47. Paik S, Shak S, Tang G, Kim C, Baker J, Cronin M, Baehner FL, Walker MG, Watson D, Park T, et al: **A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer.** *N Engl J Med* 2004, **351**:2817-2826.
48. Wechsler ME: **Managing asthma in primary care: putting new guideline recommendations into context.** *Mayo Clin Proc* 2009, **84**:707-717.
49. **Physician Fee Schedule Search.** *Centers for Medicare & Medicaid Services*, available at <https://www.cms.gov/apps/physician-fee-schedule/search/search-criteria.aspx> and accessed on 1/30/2017 2016.
50. Goodwin S, McPherson JD, McCombie WR: **Coming of age: ten years of next-generation sequencing technologies.** *Nat Rev Genet* 2016, **17**:333-351.
51. **Asthma in the US.** *Centers for Disease Control and Prevention Vital signs* <http://www.cdc.gov/vitalsigns/asthma/>, downloaded 1/30/2017 2011.
52. Jain VV, Allison DR, Andrews S, Mejia J, Mills PK, Peterson MW: **Misdiagnosis Among Frequent Exacerbators of Clinically Diagnosed Asthma and COPD in Absence of Confirmation of Airflow Obstruction.** *Lung* 2015, **193**:505-512.
53. Brower V: **Biomarkers: Portents of malignancy.** *Nature* 2011, **471**:S19-21.
54. Muraro A, Lemanske RF, Jr., Hellings PW, Akdis CA, Bieber T, Casale TB, Jutel M, Ong PY, Poulsen LK, Schmid-Grendelmeier P, et al: **Precision medicine in patients with allergic diseases: Airway diseases and atopic dermatitis-PRACTALL document of the European Academy of Allergy and Clinical Immunology and the American Academy of Allergy, Asthma & Immunology.** *J Allergy Clin Immunol* 2016, **137**:1347-1358.

55. Himes BE, Hunninghake GM, Baurley JW, Rafaels NM, Sleiman P, Strachan DP, Wilk JB, Willis-Owen SA, Klanderma B, Lasky-Su J, et al: **Genome-wide association analysis identifies PDE4D as an asthma-susceptibility gene.** *Am J Hum Genet* 2009, **84**:581-593.
56. Fromer M, Roussos P, Sieberts SK, Johnson JS, Kavanagh DH, Perumal TM, Ruderfer DM, Oh EC, Topol A, Shah HR, et al: **Gene expression elucidates functional impact of polygenic risk for schizophrenia.** *Nat Neurosci* 2016.
57. Langmead B, Trapnell C, Pop M, Salzberg SL: **Ultrafast and memory-efficient alignment of short DNA sequences to the human genome.** *Genome Biol* 2009, **10**:R25.
58. Trapnell C, Pachter L, Salzberg SL: **TopHat: discovering splice junctions with RNA-Seq.** *Bioinformatics* 2009, **25**:1105-1111.
59. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter L: **Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation.** *Nat Biotechnol* 2010, **28**:511-515.
60. DeLuca DS, Levin JZ, Sivachenko A, Fennell T, Nazaire MD, Williams C, Reich M, Winckler W, Getz G: **RNA-SeQC: RNA-seq metrics for quality control and process optimization.** *Bioinformatics* 2012, **28**:1530-1532.
61. Pedregosa F, Varoquaux Ge, Gramfort A, Michel V, Thirion B, others: **Scikit-learn: Machine Learning in Python.** *Journal of Machine Learning Research* 2011, **12**:2825-2830.
62. Guyon I, Weston, J, Barnhill, S, Vapnik, V: **Gene selection for cancer classification using support vector machines.** *Machine Learning* 2002, **46**:389-422.
63. Bewick V, Cheek L, Ball J: **Statistics review 14: Logistic regression.** *Crit Care* 2005, **9**:112-118.

64. Burges CJ: **A tutorial on support vector machines for pattern recognition.** *Data mining and knowledge discovery* 1998, **2**:121-167.
65. Freund Y, Schapire RE: **A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting.** *J Comput Syst Sci* 1997, **55**:119-139.
66. Breiman L: **Random Forests.** *Machine Learning* 2001, **45**:5-32.
67. Hollander M, Wolfe DA, Chicken E: *Nonparametric statistical methods.* John Wiley & Sons; 2013.
68. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA, Phillippy KH, Sherman PM, Holko M, et al: **NCBI GEO: archive for functional genomics data sets--update.** *Nucleic Acids Res* 2013, **41**:D991-995.

Figure Legends

Figure 1: Study flow for the identification of a nasal brush-based classifier of asthma by machine learning analysis of RNAseq data.

Subjects with mild/moderate asthma and controls without asthma were recruited for phenotyping, nasal brushing, and RNA sequencing of nasal brushings. The RNAseq data generated were then *a priori* split into development and test sets. The development set was used for differential expression analysis and machine learning (involving feature selection, classification, and statistical analyses of classification performance) to identify an asthma gene panel that can accurately classify asthma from no asthma. The asthma gene panel was then tested on eight validation test sets, including (1) the RNAseq test set of independent subjects with and without asthma, (2) two external test sets of subjects with and without asthma with nasal gene expression profiled by microarray, and (3) five external test sets of subjects with non-asthma respiratory conditions (allergic rhinitis, upper respiratory infection, cystic fibrosis, and smoking) and nasal gene expression profiled by microarray.

Figure 2: Receiver operating characteristic (ROC) curve of the predictions generated by applying the asthma gene panel to the samples in the RNAseq test set of independent subjects (n=40).

The ROC curve for a random model is shown for reference. The curve and its corresponding AUC score show that the panel performs well for both asthma and no asthma (control) samples in this test set.

Figure 3: Validation of the asthma gene panel on test sets of independent subjects with asthma.

Performance of the asthma panel in classifying asthma (A) and no asthma (C) in terms of F-measure, a conservative mean of precision and sensitivity. F-measure ranges from 0 to 1, with higher values indicating superior classification performance. The panel was applied to an RNAseq test set of independent subjects with and without asthma, and two external microarray

data sets from subjects with and without asthma (Asthma1 and Asthma2). Positive (B) and negative (D) predictive values are also provided for context.

Figure 4: Comparative performance of the asthma gene panel and other classification

models in the RNAseq test set. Performances of the asthma gene panel and other classification models in classifying asthma (left panel) and no asthma (right panel) are shown in terms of F-measure, with individual measures shown in the bars. The number of genes in each model is shown in parentheses within the bars. The asthma gene panel is labeled in red and classification models learned from the machine learning pipeline using other combinations of feature selection and classification are labeled in black. These other classification models were combinations of two feature selection algorithms (LR-RFE and SVM-RFE) and four global classification algorithms (Logistic Regression, SVM-Linear, AdaBoost and Random Forest). For context, alternative classification models (labeled in blue) are also shown and include: (1) a model derived from an alternative, single-step classification approach (sparse classification model learned using the L1-Logistic regression algorithm), and (2) models substituting feature selection with each of the following preselected gene sets - all genes after filtering, all differentially expressed genes in the development set, and known asthma genes [28] - with their respective best performing global classification algorithms. These results show the superior performance of the asthma gene panel compared to all other models, in terms of classification performance and model parsimony (number of genes included). LR = Logistic Regression. SVM = Support Vector Machine. RFE = Recursive Feature Elimination. RF = Random Forest.

Figure 5: Validation of the asthma gene panel on test sets of independent subjects with

non-asthma respiratory conditions. Performance statistics of the panel when applied to external microarray-generated data sets of nasal gene expression derived from case/control cohorts with non-asthma respiratory conditions. Performance is shown in terms of F measure (A

and C), a conservative mean of precision and sensitivity, as well as positive (B) and negative predictive values (D). The panel had a low to zero rate of misclassifying other respiratory conditions as asthma, supporting that the panel is specific to asthma and would not misclassify other respiratory conditions as asthma.

Figure 6: Heatmap showing expression profiles of the 90 gene members of the asthma gene panel. Columns shaded pink at the top denote asthma samples, while samples from subjects without asthma are denoted by columns shaded grey. 22 and 24 of these genes were over- and under-expressed in asthma samples (DESeq2 FDR ≤ 0.05), denoted by orange and purple groups of rows, respectively. The four genes in this set that have been previously associated with asthma [28] are marked in blue. The panel's inclusion of genes not previously known to be associated with asthma as well as genes not differentially expressed in asthma (beige group of rows) demonstrates the ability of our machine learning methodology to move beyond traditional analyses of differential expression and current domain knowledge.

Table 1: Baseline characteristics of subjects in the RNAseq development and test sets

	Development set			Test Set			Development vs. Test Set P value ^B
	All (n=150)	Asthma (n=53)	No Asthma (n=97)	All (n=40)	Asthma (n=13)	No Asthma (n=27)	
Age: years	26.9 (5.4)	25.7 (2.0)	27.6 (6.5)	26.2 (5.1)	25.3 (2.1)	26.6 (6.1)	0.47
Sex: female	89 (59.3%)	24 (45.3%)	65 (67.0%)	21 (52.5%)	2 (15.3%)	19 (70.4%)	0.40
Race							0.60
Caucasian	116 (77.3%)	21 (40.4%)	96 (99.0%)	32 (80.0%)	5 (38.5%)	27 (100.0%)	
African American	24 (16.0%)	23 (43.4%)	1 (1.0%)	5 (12.5%)	5 (38.5%)	0 (0.0%)	
Latino	5 (3.3%)	5 (9.4%)	0 (0.0%)	3 (7.5%)	3 (23.1%)	0 (0.0%)	
Other	5 (3.3%)	4 (7.5%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	
FEV1 ^A : % predicted	94.7 (10.0)	94.6% (10.9)	94.8 (9.7)	94.5 (11.4)	94.4 (12.0)	94.6 (11.3)	0.90
FEV1/FVC ^A : %	82.5 (6.4)	81.5 (6.7)	83.1 (6.3)	82.7 (5.5)	84.8 (4.4)	81.6 (5.8)	0.91
Bronchodilator response: %	5.6 (6.0)	8.7 (6.4)	3.9 (5.1)	4.5 (5.4)	7.0 (6.1)	3.3 (4.7)	0.29
Age asthma onset: years		3.2 (2.7)	n/a		3.4 (2.0)		0.78
Allergic rhinitis	60 (40.0%)	29 (54.7%)	31 (32.0%)	7 (17.5%)	7 (53.8%)	0 (0.0%)	0.009
Nasal steroids	14 (9.3%)	9 (17.0%)	5 (5.2%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	0.07
Smoking	7 (4.7%)	1 (1.9%)	6 (6.2%)	1 (2.5%)	0 (0.0%)	1 (3.7%)	1.0

Mean (SD) or Number (%) provided

^Apre-bronchodilator measures. FEV1 = forced expiratory flow volume in 1 second, FVC = forced vital capacity

^BFisher's Exact test for categorical variables and t-test for continuous variables

Supplementary Materials

Supplementary Figure 1: variancePartition analysis of the RNAseq development set.

Supplementary Figure 2: Visual description of the machine learning pipeline used to select predictive features (genes) and develop classification models based on them in the RNAseq development set.

Supplementary Figure 3: Visual description of the feature (gene) selection component of the machine learning pipeline.

Supplementary Figure 4: Critical Difference plots demonstrating results of the statistical comparison of the performance of 100 asthma classification models obtained by various combinations of feature selection and global classification algorithms in terms of the classification performance and parsimony (numbers of genes included) of the models.

Supplementary Figure 5: Evaluation measures for classification models.

Supplementary Figure 6: Performance of permutation-based random classification models in test sets of independent subjects with asthma and controls.

Supplementary Figure 7: Performance of permutation-based random classification models in test sets of independent subjects with non-asthma respiratory conditions and controls.

Supplementary Figure 8: Distribution of DESeq2 FDR values of differential expression in the asthma gene panel (blue bars) vs. other genes in the RNAseq development set (coral bars).

Supplementary Table 1: Lists of over- and under-expressed genes and pathways in asthma cases compared to controls (in different tabs of this file). Differentially expressed genes were identified using DESeq2 [24] applied to the development set, and enriched pathways were identified from the Molecular Signature Database [25], both using an upper FDR threshold of 0.05.

Supplementary Table 2: List of known asthma-associated genes from a recent review of asthma genetics [28] that overlap with genes in our RNAseq data sets.

Supplementary Table 3: Characteristics of the external asthma cohorts used in the validation of the asthma gene panel.

Supplementary Table 4: Characteristics of the external cohorts with non-asthma respiratory conditions and controls used in the validation of the asthma gene panel.











