## ORIGINAL ARTICLE

# Genetic regulation of gene expression in the lung identifies CST3 and CD22 as potential causal genes for airflow obstruction

Maxime Lamontagne,<sup>1</sup> Wim Timens,<sup>2</sup> Ke Hao,<sup>3</sup> Yohan Bossé,<sup>1,4</sup> Michel Laviolette,<sup>1</sup> Katrina Steiling,<sup>5</sup> Joshua D Campbell,<sup>5</sup> Christian Couture,<sup>1</sup> Massimo Conti,<sup>1</sup> Karen Sherwood,<sup>6</sup> James C Hogg,<sup>6,7</sup> Corry-Anke Brandsma,<sup>2</sup> Maarten van den Berge,<sup>8</sup> Andrew Sandford,<sup>6,9</sup> Stephen Lam,<sup>10</sup> Marc E Lenburg,<sup>5</sup> Avrum Spira,<sup>5</sup> Peter D Paré,<sup>6,9</sup> David Nickle,<sup>11</sup> Don D Sin,<sup>6,9</sup> Dirkje S Postma<sup>8</sup>

## ABSTRACT

Additional material is published online only. To view please visit the journal online (http://dx.doi.org/10.1136/ thoraxjnl-2014-205630).

For numbered affiliations see end of article.

### Correspondence to

Professor Dirkie S Postma. University of Groningen, University Medical Center Groningen, Department of Pulmonary Medicine and Tuberculosis, GRIAC Research Institute, AA11, Hanzeplein, PO Box 30001, Groningen 9700 RB, The Netherlands; d.s.postma@umcq.nl

ML, WT, and KH contributed equally.

Received 28 April 2014 Revised 16 July 2014 Accepted 14 August 2014 Published Online First 2 September 2014

**Background** COPD is a complex chronic disease with poorly understood pathogenesis. Integrative genomic approaches have the potential to elucidate the biological networks underlying COPD and lung function. We recently combined genome-wide genotyping and gene expression in 1111 human lung specimens to map expression quantitative trait loci (eOTL).

**Objective** To determine causal associations between COPD and lung function-associated single nucleotide polymorphisms (SNPs) and lung tissue gene expression changes in our lung eQTL dataset.

Methods We evaluated causality between SNPs and gene expression for three COPD phenotypes: FEV<sub>1</sub>% predicted, FEV<sub>1</sub>/FVC and COPD as a categorical variable. Different models were assessed in the three cohorts independently and in a meta-analysis. SNPs associated with a COPD phenotype and gene expression were subjected to causal pathway modelling and manual curation. In silico analyses evaluated functional enrichment of biological pathways among newly identified causal genes. Biologically relevant causal genes were validated in two separate gene expression datasets of lung tissues and bronchial airway brushings. **Results** High reliability causal relations were found in SNP-mRNA-phenotype triplets for FEV<sub>1</sub>% predicted (n=169) and FEV<sub>1</sub>/FVC (n=80). Several genes of potential biological relevance for COPD were revealed. eQTL-SNPs upregulating cystatin C (CST3) and CD22 were associated with worse lung function. Signalling pathways enriched with causal genes included xenobiotic metabolism, apoptosis, protease-antiprotease and oxidant-antioxidant balance.

**Conclusions** By using integrative genomics and analysing the relationships of COPD phenotypes with SNPs and gene expression in lung tissue, we identified CST3 and CD22 as potential causal genes for airflow obstruction. This study also augmented the understanding of previously described COPD pathways.



To cite: Lamontagne M, Timens W, Hao K, et al. Thorax 2014;69:997-1004.

## INTRODUCTION

Genome-wide association studies (GWAS) have revolutionised our ability to identify common genetic variants that are associated with complex chronic diseases.<sup>1</sup> This approach has been applied

## **Key messages**

## What is the key question?

What are the causal genetic variants changing gene expression in the lung that in turn associate with lower lung function and COPD?

## What is the bottom line?

► Lung function-associated genetic variants alter the mRNA expression of nearby genes involved in biological pathways underpinning pulmonary function and COPD pathogenesis.

## Why read on?

New genes of airflow obstruction are identified with a generalised framework for the identification of causal genes from joint examination of genome-wide genotyping and gene expression data in the same patients.

to COPD, a lung disease that is caused predominantly by cigarette smoking in the western world.<sup>2-</sup>

<sup>4</sup> It is well known that only a subset of heavy smokers (15–20%) develop clinically relevant COPD and there is considerable evidence that there is a substantial genetic component involved in its pathogenesis.<sup>5</sup> Although GWAS have identified novel loci that harbour susceptibility genes, they do not allow precise identification of the causal variant (or variants). In addition, GWAS do not provide information on how and to what extent the gene (or genes) within the susceptibility loci contribute to the phenotype. Interestingly, the majority of genetic variants which have been associated with disease traits by GWAS do not affect the coding sequence of genes but are located in intergenic regions or introns.<sup>6</sup> Possible explanations are that the associated alleles are in linkage disequilibrium (LD) with rarer coding alleles with large effect sizes<sup>7</sup> and/or that the genetic variants control the level of expression of genes involved in pathogenetic pathways. For complex genetic diseases, such as COPD, the effects of susceptibility alleles may primarily act by regulating gene





expression rather than by altering protein coding as in most Mendelian diseases.  $^{\rm 8}$ 

We recently reported the discovery of a large number of lungspecific expression quantitative trait loci (eQTLs)<sup>9</sup> and identified the most likely causal genes within three GWAS-nominated COPD susceptibility loci.<sup>10</sup> The aim of the present study is to use the power of genome-wide mRNA expression arrays combined with genome-wide interrogation of single nucleotide polymorphisms (SNPs) to pinpoint specific SNPs that are related to lung tissue gene expression and to COPD phenotypes. Identification of susceptibility alleles that function as strong eQTLs increases the likelihood of identifying the true susceptibility gene within loci with large areas of LD.<sup>8</sup> Moreover, the use of integrative genomics by combining environmental exposure data with susceptibility alleles, RNA expression levels, and different COPD phenotypes can unravel causal genetic relationships.<sup>11</sup> The basic principle of the current study is to map the genetic regulation of gene expression to identify DNA variants that induce changes in transcriptional networks that in turn contribute to COPD and airway obstruction pathogenesis.

### **METHODS**

### Subject selection

The methods for subject selection and phenotyping, and for interrogation of gene expression and genotype were recently described.<sup>9</sup> The lung tissue used for discovery of eQTLs was from 1111 human subjects who underwent lung surgery at three academic sites, Laval University, University of British Columbia (UBC) and University of Groningen, henceforth referred to as Laval, UBC and Groningen, respectively. All lung specimens from Laval were obtained from patients undergoing lung cancer surgery and were harvested from a site distant from the tumour. At UBC, the majority of samples were from patients undergoing resection of small peripheral lung lesions. Additional samples were from autopsy and at the time of lung transplantation. At Groningen, the lung specimens were obtained at surgery from patients with various lung diseases, including patients undergoing therapeutic resection for lung tumours, harvested from a site distant from the tumour, and lung transplantation. For the present study, the principal aim was to examine smoking-related airway obstruction. Thus, we excluded subjects whose lung function may have been influenced by lung diseases other than COPD and lung cancer. Exclusion criteria are provided in the online supplement.

### **COPD** phenotypes and GWAS

Genome-wide association was performed using linear or logistic regression models on the three phenotypes:  $FEV_1\%$  predicted and  $FEV_1/FVC$  as continuous variables, and COPD defined dichotomously based on an  $FEV_1/FVC<0.7$  cut-off (see online supplement). Single-marker association tests were run within each cohort adjusting for age, gender and smoking status. Fixed-effects meta-analysis was then performed combining the three cohorts using inverse SE weighting.

### **Expression trait processing**

Expression traits were adjusted for age, gender and smoking status as described previously.<sup>9</sup> Gene expression data are available in the Gene Expression Omnibus repository through accession number GSE23546.

### **Causality models**

We evaluated three competing causality models to describe the relationships between lung eQTL-SNPs, RNA expression and



**Figure 1** Causality models showing the relationships between the expression of a gene, a phenotype and a single nucleotide polymorphism (SNP). Three models are depicted including a causal model (M1), a reactive model (M2), and an independent model (M3). M1 is the simplest, and states that the genotype at an expression quantitative trait loci SNP acts directly on gene expression pattern to produce the phenotype. M2 states that the gene expression pattern is reactive to the phenotype, and M3 states that the gene expression pattern and phenotype are independent.

COPD phenotypes<sup>12</sup> (figure 1). Model 1 indicates a potential causal relationship where the SNP acts on gene expression to produce the phenotype, which was the main interest of this work. Model 2 indicates that the gene expression pattern is reactive to the phenotype, that is, the phenotype drives gene expression. Model 3 is the independent model where the SNP acts on the phenotype and gene expression independently. We required a p value  $<1 \times 10^{-3}$  for SNP associated with phenotype before examining a particular variable triplet (ie, SNP, gene expression and phenotypes). The causal model (model 1) was selected when p value <0.05 was found for SNP associated with gene expression adjusting for phenotype and p value >0.05 was found for SNP associated with phenotype adjusting for gene expression. More details are provided in the online supplement. Analyses were performed in the three cohorts separately and then combined into a meta-analysis. We then conducted sample bootstrapping and repeated the causality test for N=1000 realisations. The reliability score of the molecular relationship is the fraction of bootstrap realisations that supports the call observed. Threshold for p values, number of bootstrap reselections and reliability cut-off (>0.8) were selected based on previous literature.<sup>12</sup>

## Manual curation and pathway analyses on reliable causal models

The biology and possible role of genes in the causal model were reviewed by manual curation and bioinformatics tools, including Ingenuity Pathway Analysis (IPA), MetaCore and Partek (see online data supplement).

### **Replication studies**

Biologically relevant causal genes were validated in two datasets. First, a bronchial airway epithelial dataset where genome-wide gene expression levels were obtained from bronchial brushing of 238 individuals and associated with lung function as previously described.<sup>13</sup> Second, a regional lung tissue dataset where genome-wide gene expression was obtained for eight regions of the same lung and for eight patients. The association between gene expression levels and micro-CT-based regional emphysema severity was assessed.<sup>14</sup> More details are provided in the online data supplement.

### RESULTS

Of the 1111 subjects in whom there were data on genotype and gene expression, 848 with sufficient phenotypic information

orax
: firs
st pc
ublis
hed
as
10.1
1136
ò∕thc
orax
jn-2
2014
<b>1-</b> 20
563
õ
n 2
Sep
tem
iber
201
4. [
Dow
nloa
idec
d fro
m <mark>h</mark>
ttp:/
/tho
rax.
bmj.
8
<b>n</b> / 0
n A
orii 2
27, 2
202
4 by
gue
est.
Pro
tect
ed b
уу с
yqc

H

able 1 Clinical characte	eristics of patients
--------------------------	----------------------

Characteristics	Laval (n=403)	UBC (n=270)	Groningen (n=175)
Age (years)	63.4±9.8	63.9±10.0	59.7±10.0
Male/female (n)	224/179	144/126	83/92
Body mass index (kg/m <sup>2</sup> )	26.6±5.2	25.7±5.4	24.6±4.1
FEV <sub>1</sub> % predicted	82.2±17.6	83.6±22.3	72.4±24.9
FEV <sub>1</sub> /FVC	0.68±0.10	0.69±0.13	0.64±0.16
COPD (%)	209 (51.9)	114 (42.2)	112 (64.0)
Stage 1: mild	80 (38.2)	43 (37.7)	16 (14.3)
Stage 2: moderate	117 (56.0)	60 (52.6)	30 (26.8)
Stage 3: severe	11 (5.3)	2 (1.8)	11 (9.8)
Stage 4: very severe	1 (0.5)	9 (7.9)	44 (39.3)
Missing data	0	0	11 (9.8)
Non-COPD (%)	194 (48.1)	156 (57.8)	63 (36.0)
Smoking (%)			
Smoker	89 (22.1)	91 (33.7)	12 (6.9)
Ex-smoker	281 (69.7)	150 (55.6)	117 (66.9)
Non-smoker	33 (8.2)	16 (5.9)	43 (24.6)
Missing data	0	13 (4.8)	3 (1.7)
Pack years (years)	48.6±27.4	46.0±28.6	32.1±18.1

Continuous variables are presented as means±SDs.

UBC, University of British Columbia.

were included in the analysis. The demographic and clinical features of the subjects in the three cohorts are described in table 1.

#### Causal pathways that fit model 1

Model fitting was restricted to SNPs that were significantly associated with at least one of the lung function variables or COPD as a categorical variable at a p value cut-off of  $10^{-3}$ . SNPs were considered an eQTL if they had a p value  $\leq 10^{-5}$ . Using these

(hronic o	hetriictiva	nulmonary	I dicaaci
	$\mathbf{D}$	punnonar	y unseus

criteria, there were 4465 SNP-mRNA-phenotype triplets. Of these, 249 triplets showed a significant fit to the causal model for at least one of the phenotypes with a reliability score >0.8 in at least one cohort and/or in the meta-analysis. A fit to this model means that the SNP was associated with one of the phenotypes and affected gene expression in a direction that supported a causal relationship.

The list of causal models with a confidence score >0.8 in any one of the cohorts or in the meta-analysis is shown in online supplementary table S1. For FEV<sub>1</sub>% predicted, 169 causal models were found, including 122 unique SNPs and 169 probe sets. The 169 models included 168 *cis* eQTLs and 1 *trans* eQTL, respectively. For FEV<sub>1</sub>/FVC, 80 causal models were found involving 63 unique SNPs and 80 probe sets. Among these 80 models, 79 were *cis* eQTLs and 1 was a *trans* eQTL. No causal model was found when using COPD as a categorical variable. There was no overlap between causal pathway models for FEV<sub>1</sub>% predicted and FEV<sub>1</sub>/FVC.

### Manual curation of significant and reliable causal models

Significant causal models were inspected manually. We identified a number of models of potential biological relevance for FEV<sub>1</sub>% predicted and FEV<sub>1</sub>/FVC. The most biologically relevant models are listed in table 2. The p values and direction of effect for each of the SNP associations with the gene expression and the phenotype are also indicated. For example, SNP rs6048956 was significantly associated with cystatin C (*CST3*) transcript. The common allele was associated with higher expression of the transcript (figure 2 and positive eQTL Z score in table 2) and with lower FEV<sub>1</sub>% predicted (figure 2 and negative Z score with phenotype in table 2). Taken together, these results suggest that the common allele confers susceptibility to a lower FEV<sub>1</sub>% predicted value through upregulation of the *CST3* mRNA expression levels in the lung. This was confirmed in a second causality model that interrogated a different probe

Table 2 Models of	of biological	relevance identified	by	manua	curation
-------------------	---------------	----------------------	----	-------	----------

SNPs	Reference allele (freq)*	Gene (probe set)	Lung eQTL p value <sup>†</sup>	eQTL Z score <sup>‡</sup>	p Value phenotype§	Z score phenotype
FEV <sub>1</sub> % predicted	l					
rs769178	G (0.91)	NCR3 (100125842_TGI_at)	8.03×10 <sup>-30</sup>	11.3	4.7×10 <sup>-4</sup>	3.5
rs6048956	C (0.77–0.79)	<i>CST3</i> (100307577_TGI_at)	1.54×10 <sup>-68</sup>	17.5	7.8×10 <sup>-4</sup>	-3.4
rs6515375	G (0.79–0.81)	CST3 (100125967_TGI_at)	3.65×10 <sup>-6</sup>	4.6	5.8×10 <sup>-4</sup>	-3.4
rs2270859	G (0.83–0.88)	CSTA (100148334_TGI_at)	1.74×10 <sup>-9</sup>	6.0	3.1×10 <sup>-6</sup>	-4.7
rs4550905	G (0.26–0.31)	PPARGC1A (100131093_TGI_at)	1.69×10 <sup>-5</sup>	-4.3	4.4×10 <sup>-4</sup>	-3.5
rs3803761	G (0.71–0.77)	FLCN (100135396_TGI_at)	3.77×10 <sup>-29</sup>	11.2	3.4×10 <sup>-2</sup>	-2.1
rs1543438	A (0.76–0.82)	BCL2L1 (100158784_TGI_at)	3.89×10 <sup>-5</sup>	-4.1	7.94×10 <sup>-5</sup>	3.9
rs9880397	G (0.62–0.64)	CADM2 (100162763_TGI_at)	1.3×10 <sup>-8</sup>	-5.7	2.9×10 <sup>-4</sup>	3.6
rs2466183	T (0.84–0.86)	TNFRSF10B (100153254_TGI_at)	3.76×10 <sup>-5</sup>	-4.1	6.6×10 <sup>-4</sup>	3.4
FEV <sub>1</sub> /FVC						
rs17754977	A (0.30–0.32)	GSTO2 (100132911_TGI_at)	3.96×10 <sup>-5</sup>	4.1	8.2×10 <sup>-4</sup>	3.3
rs9987135	T (0.28–0.31)	DEPDC6 (100154484_TGI_at)	6.13×10 <sup>-59</sup>	16.2	1.1×10 <sup>-4</sup>	-3.9
rs10411704	T (0.79–0.82)	CD22 (100154732_TGI_at)	1.70×10 <sup>-40</sup>	-13.3	4.6×10 <sup>-4</sup>	3.5
rs12179536	A (0.80–0.84)	MUC22 (100304000_TGI_at)	3.47×10 <sup>-33</sup>	-12.0	3.0×10 <sup>-3</sup>	3.0
rs2287765	T (0.91–0.93)	SPINK5 (100305138_TGI_at)	2.85×10 <sup>-10</sup>	6.3	9.2×10 <sup>-4</sup>	-3.3

p Values and Z scores in this table are from the meta-analysis. As indicated in the text, the significant causal models were selected based on results of both individual cohorts and meta-analysis (SNP associated with phenotype with p value <10<sup>-3</sup>, SNPs associated with gene expression with p value  $\leq 10^{-5}$ , and reliability score >0.8 in at least one cohort and/or in the meta-analysis). Known role of genes are provided in online supplementary table S2.

\*Frequency of the reference allele in the three cohorts.

†Lung eQTL p-value from the meta-analysis.

‡Z score from the eQTL meta-analysis showing the direction of effect for the SNP on gene expression.

§p Value for association between the SNP and phenotype from the meta-analysis.

IEffect size (Z score) for association between the SNP and the phenotype from the meta-analysis, showing the direction of effect for the SNP on the phenotype.

eQTL, expression quantitative trait loci; SNP, single nucleotide polymorphism.

right



**Figure 2** Direction of effects for causality model with rs6048956, mRNA expression of cystatin C (*CST3*) and FEV<sub>1</sub>% predicted. The left, centre and right panels show the results for Laval, University of British Columbia (UBC) and Groningen samples, respectively. (A–C) are boxplots of gene expression levels in the lung for *CST3* according to genotype groups for single nucleotide polymorphism (SNP) rs6048956. The left y-axis shows the mRNA expression levels for *CST3*. The x-axis represents the three genotype groups for SNP rs6048956. The right y-axis shows the proportion of the gene expression variance explained by the SNP (black bar). Box boundaries, whiskers and centre mark in boxplots represent the first and third quartiles, the most extreme data point which is no more than 1.5 times the IQR, and median, respectively. (D–F) are barplots of the mean and SE of FEV<sub>1</sub>% predicted according to genotype groups for SNP rs6048956.

set for *CST3*, FEV<sub>1</sub>% predicted and rs6515375 (table 2 and see online supplementary figure S1). Analogous observations were made for natural cytotoxicity triggering receptor 3 (*NCR3*), cystatin A (*CSTA*), peroxisome proliferation-activated receptor  $\gamma$ coactivator 1 $\alpha$  (*PPARGC1A*), folliculin (*FLCN*), BCL2-like 1 (*BCL2L1*), cell adhesion molecule 2 (*CADM2*), and tumour necrosis factor receptor superfamily member 10b (*TNFRSF10B*) for FEV<sub>1</sub>% predicted, and glutathione S-transferase  $\omega$  2 (*GSTO2*), DEP domain containing 6 (*DEPDC6*), CD22 molecule (*CD22*), mucin 22 (*MUC22*), and serine peptidase inhibitor Kazal type 5 (*SPINK5*) for FEV<sub>1</sub>/FVC (table 2). The direction of effects for all these causality models are illustrated in online supplementary figures S2–13.

### **Pathway analyses**

To find key biological pathways involved in COPD, causality genes were overlaid on canonical pathways available in IPA. One hundred and sixty-nine genes and 80 causality genes for FEV<sub>1</sub> and FEV<sub>1</sub>/FVC were considered, respectively (see online supplementary table S1). Table 3 shows all canonical pathways enriched with causality genes (p<0.05). Of interest was the aryl

hydrocarbon receptor signalling pathway, which is involved in xenobiotic clearance. Five causality genes were noted in this canonical pathway, including *ARNT* (aryl hydrocarbon receptor nuclear translocator), *IL6* (interleukin 6), *GSTO2*, *ALDH8A1* (aldehyde dehydrogenase 8 family, member A1) and *TRIP11* (thyroid hormone receptor interactor 11). The aryl hydrocarbon receptor signalling pathway and the location of the causality genes are illustrated in online supplementary figure S14. Another pathway involved in xenobiotic handling, the xenobiotic metabolism signalling pathway, was also enriched for causality genes, many of which overlapped with those in the aryl hydrocarbon receptor signalling pathway (*ARNT*, *IL6*, *GSTO2* and *ALDH8A1*). The former also includes *PPARGC1A* and *PRKCE* (protein kinase C,  $\varepsilon$ ). This pathway is illustrated in online supplementary figure S15.

For the Partek analyses 118 and 45 transcripts were identified for FEV<sub>1</sub>% predicted and FEV<sub>1</sub>/FVC, respectively. These transcripts mapped to a total of 115 genes in Partek GS (based on official gene symbols) and were parsed to Partek Pathway Suite. Causality genes were overlaid onto canonical pathways from the REACTOME and KEGG databases. Table 4 shows the canonical

p Value 0.00002 0.00015 0.00069 0.00151 0.00468 0.00589 0.00617 0.00661 0.00676 Communication between innate and adaptive immune cells 0.00676 0.00759 0.00891 0.00977 0.01023 0.01047 0.01072 0.01230 0.01230 0.01995 0.02455 0.02692 0.02951 0.03162 0.03631 0.04365

0.04467

0.04677

Pathways discussed in the text are given in italic.

pathways enriched with causality genes (p<0.05). TNF-related apoptosis-inducing ligand (TRAIL) signalling was the most significant pathway, with several other apoptosis-related pathways included in the list of top significant associations. The glutathione metabolism pathway connects to both the cyanoamino acid and taurine and hypotaurine metabolism pathways (see online supplementary figure S16). These three pathways were identified in our enrichment analysis and were also identified using the IPA system. GO functional categories also showed enrichment for causality genes (see online supplementary figure S17) including the  $\gamma$ -glutamyl transferase activity.

 Table 3
 Canonical pathways enriched for causality genes using

the Ingenuity Pathway Analysis system

Graft-versus-host disease signalling

Taurine and hypotaurine metabolism

Aryl hydrocarbon receptor signalling

Amyotrophic lateral sclerosis signalling

Systemic lupus erythematosus signalling

Allograft rejection signalling

Cyanoamino acid metabolism

Macropinocytosis signalling

Selenoamino acid metabolism

Aminoacyl-tRNA biosynthesis

Actin cytoskeleton signalling

FAK signalling

Apoptosis signalling

Xenobiotic metabolism signalling

Cellular effects of sildenafil (Viagra)

Glutathione metabolism

PXR/RXR activation

Autoimmune thyroid disease signalling

Regulation of actin-based motility by Rho

OX40 signalling pathway

Antigen presentation pathway

Crosstalk between dendritic cells and natural killer cells

Neuroprotective role of THOP1 in Alzheimer's disease

Cytotoxic T-lymphocyte-mediated apoptosis of target cells

Pathwavs

Cdc42 signalling

Integrin signalling

## Replication of the most biologically relevant causality models

No other large-scale lung eQTL dataset is available in patients with and without COPD. To replicate the most biologically relevant causality models, we relied on two genome-wide expression datasets. Causality genes were first validated in a gene expression study of bronchial airway epithelial cells obtained by bronchoscopy. All 13 genes represented in table 2, except MUC22, were assayed in this airway dataset. Two genes were significantly associated with lung function at a false discovery rate (FDR) of 5% in the bronchial airway epithelial dataset: CSTA (FDR=0.002 with FEV<sub>1</sub>% predicted) and TNFRSF10B (FDR= $9.35 \times 10^{-5}$  with FEV<sub>1</sub>% predicted). In both cases, the direction of effect was the same as that in the current study (table 5). For example, higher CSTA mRNA levels were

#### Table 4 Canonical pathways enriched for causality genes using Partek

Pathways	p Value	Pathway ID
TRAIL signalling	0.0003	reactome_pathway_506
Cyanoamino acid metabolism	0.001	kegg_pathway_36
Cell adhesion molecules	0.002	kegg_pathway_139
Extrinsic pathway for apoptosis	0.002	reactome_pathway_502
Death receptor signalling	0.002	reactome_pathway_503
Taurine and hypotaurine metabolism	0.002	kegg_pathway_34
Glutathione metabolism	0.003	kegg_pathway_39
Aminoacyl-tRNA biosynthesis	0.004	kegg_pathway_81
Apoptosis	0.006	reactome_pathway_501
Cytosolic tRNA aminoacylation	0.006	reactome_pathway_1127
Glutathione conjugation	0.007	reactome_pathway_256
Natural killer cell mediated cytotoxicity	0.011	kegg_pathway_152
The NLRP1 inflammasome	0.015	reactome_pathway_469
tRNA aminoacylation	0.018	reactome_pathway_1126
Apoptosis	0.023	kegg_pathway_126
BoNT light chain types B, D and F cleave VAMP/synaptobrevin	0.024	reactome_pathway_647
Amyloids	0.026	reactome_pathway_643
Nucleotide-binding domain, leucine rich repeat containing receptor (NLR) signalling pathways	0.026	reactome_pathway_465
BH3-only proteins associate with and inactivate anti-apoptotic BCL-2 members	0.034	reactome_pathway_517
Axonal growth inhibition (RHOA activation)	0.039	reactome_pathway_984
Vitamin C (ascorbate) metabolism	0.039	reactome_pathway_188
Import of palmitoyl-CoA into the mitochondrial matrix	0.039	reactome_pathway_70
Downregulation of ERBB4 signalling	0.039	reactome_pathway_1023
p75NTR regulates axonogenesis	0.043	reactome_pathway_982
Caspase-8 is formed from procaspase-8	0.043	reactome_pathway_507
Activation of procaspase-8	0.043	reactome_pathway_508
Endosomal/vacuolar pathway	0.043	reactome_pathway_413
Phase II conjugation	0.047	reactome_pathway_249
SLBP independent processing of histone pre-mRNAs	0.048	reactome_pathway_1111
Arachidonic acid metabolism	0.050	kegg_pathway_56

Pathways in italic are also found in ingenuity pathway analysis. TRAIL, tumour necrosis factor related apoptosis-inducing ligand.

associated with worse lung function (table 2 and see online supplementary figure S3).

Causality genes identified in this study were also compared with a second lung transcriptomic study that evaluated the impact of regional emphysema severity on gene expression. Interestingly, CD22 was positively associated with regional emphysema severity within individuals ( $p=5.8 \times 10^{-5}$ ) and was a part of the 127 gene signature for emphysema identified in that study.<sup>14</sup> This observation is consistent with the current study showing that carriers of the rare allele for rs10411704 had greater mRNA expression of CD22 and worse FEV<sub>1</sub>/FVC (table 2 and see online supplementary figure S11). Associations with regional emphysema severity were also observed for three other genes in table 2 including NCR3 (p=0.058), PPARGC1A (p=0.081) and BCL2L1 (p=0.051), but these were not statistically significant. However, the direction of effect was consistent only for PPARGC1A. The expression of this gene was found to decrease with emphysema severity in the regional lung tissue dataset and, in the current study, carriers of the rare allele for

Table 5	Replication of causality genes in the bronchial airway
epithelial	and the regional lung tissue datasets

	Replication datasets		
	Bronchial airway epithelium <sup>13</sup>	Regional lung tissue <sup>14</sup>	
FEV <sub>1</sub> % predicted			
NCR3	True	False	
CST3	True		
CSTA	True*		
PPARGC1A	False	True	
FLCN	False		
BCL2L1	True	False	
CADM2	False		
TNFRSF10B	True*		
FEV <sub>1</sub> /FVC			
GSTO2	False		
DEPDC6	False		
CD22	True	True*	
MUC22			
SPINK5	True		

rs4550905 were associated with less mRNA expression of *PPARGC1A* in the lung and lower  $FEV_1$  (table 2 and see online supplementary figure S4).

### DISCUSSION

This investigation integrated a genome-wide eQTL study on lung tissue with genome-wide genetic association results for lung function and COPD in the same subjects. For the discovery of eQTLs we used the entire dataset which consisted of 1111 tissue samples from subjects who had lung surgery for a variety of reasons. In order to focus on COPD, we limited the study to 848 subjects with sufficient phenotypic information and without a lung disease (other than COPD and lung cancer) which could cause abnormalities of pulmonary function. We limited the causal pathway analysis to SNPs that were both significantly related to gene expression  $(p < 10^{-5})$  and were associated with one of three COPD phenotypes ( $p < 10^{-3}$ ), that is, FEV<sub>1</sub>% predicted, FEV1/FVC and COPD defined as FEV1/FVC<0.7. Of the 4465 SNP-mRNA-phenotype triplets which met our inclusion criteria, 249 triplets showed a significant fit to the causality model for at least one of the phenotypes with a reliability score  $\geq 0.8$ . We thus provide evidence that these SNPs influence disease susceptibility by altering gene expression. Causality pathway genes were enriched in pathways involved in xenobiotic handling, antiprotease and antioxidant activity and apoptosis.

Two of the eQTL-SNPs in *CST3* (rs6515375 and rs6048956) and another for *CSTA* (rs2270859) were in the causal pathway for FEV<sub>1</sub>% predicted. The two *CST3* SNPs were in perfect LD and were eQTLs for two different probe sets. *CST3* and *CSTA* are cysteine antiproteases; *CST3* antagonises cysteine cathepsins such as cathepsin L and S, and *CSTA* acts similarly for cathepsins B, H and L.<sup>15</sup> <sup>16</sup> Interestingly the direction of association is such that the alleles that are associated with a higher mRNA level of *CST3* and *CSTA* are associated with lower FEV<sub>1</sub>/FVC.

The above findings might seem paradoxical since cystatins are antiproteases and if one simply invoked the protease–antiprotease hypothesis, then one might expect that individuals with higher levels to be protected from COPD. One possibility is that

upregulation of CST3 mRNA and protein could be the result of a feedback loop stimulated by high protease levels. This would be supported by observations in bronchoalveolar lavage fluid and serum reporting higher cystatin C level in patients with emphysema.<sup>17 18</sup> Alternatively excessive inhibition of proteases may confer biological effects on a yet to be discovered mechanism which contributes to the pathogenesis of airflow obstruction. Further molecular studies of the involved proteins might shed more light on this. We previously reported SNPs associated with CST3 mRNA and protein levels in alveolar macrophages.<sup>11</sup> However, in the latter study, higher CST3 mRNA in alveolar macrophages was associated with higher FEV<sub>1</sub>. Together, these studies suggest that the cystatin genes are important in the pathogenesis of COPD, however the relationship between CST3 mRNA expression levels and lung function remains to be elucidated in relevant tissues and cell types.

With respect to antioxidant activity, eQTLs from PPARGC1A and GSTO2 were found. PPARGC1A binds to peroxisome proliferation-activated receptor  $\gamma$  (PPAR $\gamma$ ) by induction of PPAR $\gamma$ ligands, coactivating PPARy target genes, involved in antioxidant activity. Compared with controls, expression levels of PPARy, PGC-1 $\alpha$  and  $\gamma$  glutamylcysteine synthetase ( $\gamma$ -GCS) have been reported to be significantly increased in the lungs of patients with mild COPD, and progressively decreased in more severe disease.<sup>20</sup> These authors concluded that  $\gamma$ -GCS showed compensatory upregulation in the early stage of COPD, which progressively decompensated with disease progression and that the activation of the PPARy/PGC-1a pathway may protect against COPD progression by upregulating  $\gamma$ -GCS and relieving oxidative stress. Furthermore PPARGC1A has been described to be involved in bronchial smooth muscle remodelling and skeletal muscle wasting.<sup>21 22</sup> For GSTO2, a strong association of the Asn142Asp SNP with FEV1 and FVC was found in the Framingham Heart Study.<sup>23</sup> In a latter study, the Asn142Asp polymorphism in GSTO2 and the GSTO1 140Asp/ GSTO2 142Asp haplotype were associated with increased risk of COPD but failed to reveal an association between lung function parameters and non-synonymous coding SNPs in the GSTO genes.<sup>24</sup> Polymorphisms in GSTO2 were also associated with COPD either with or without lung cancer.<sup>25</sup>

Pathways involved in apoptosis and in handling of xenobiotic pathways were significantly enriched for causal genes. The results of the present study aid in understanding the underlying genetic dysregulation of apoptosis. In general, increased apoptosis of endothelial cells and fibroblasts has been shown to contribute to the development of emphysema.<sup>26</sup><sup>27</sup> This is partly due to an imbalance caused by excess oxidant and protease effects related to cigarette smoking and to intrinsic dysregulation of apoptosis induced in susceptible individuals. This may explain why the apoptotic effects are larger in smokers with COPD than smokers without COPD. Dysregulation of apoptosis can also work the other way: decreased apoptosis in cells of the immune system can lead to sustained inflammation, and when occurring in fibroblasts (eg, in the bronchial wall) can induce fibrosis. Emphysema is characterised by alveolar cell apoptosis, which was shown to be mediated by increased levels of apoptotic proteins including TRAIL receptors.<sup>28</sup> Interestingly, TRAIL signalling was the most significant pathway enriched with causality genes using Partek. TRAIL signalling has been shown to regulate immune responses in the lung and can lead to a sustained, proinflammatory response that contributes to vascular disease. An opposite effect is related to the death receptor signalling pathway, which, in humans, binds with TRAIL, induces formation of a death-inducing signalling complex, ultimately leading to caspase activation and initiation of apoptosis.<sup>29</sup>

The aryl hydrocarbon receptor signalling pathway is involved in xenobiotic clearance and therefore of relevance to the known vulnerability of patients with COPD to (cigarette) smoke. Furthermore additional metabolic pathways including the glutathione metabolism pathway and two subpathways involving metabolism of the cyanoamino acids, taurine and hypotaurine were also significantly enriched in the causal analysis. Genetic variants in a number of genes in the glutathione pathway have previously been associated with risk for COPD.<sup>25 30</sup>

No other large-scale lung eQTL dataset is available in patients with and without COPD. To replicate the most biologically relevant causality models, we relied on two genome-wide expression datasets. Causality genes were first validated in a gene expression study of bronchial airway epithelial cells obtained by bronchoscopy and then in a transcriptomic study of whole lung explanted at surgery. Whole genome genotyping is not available for these datasets. Accordingly, replication of significant triplets (ie, SNP–mRNA–phenotype) can only evaluate the concordance between gene expression and phenotype. Additional studies with phenotype, genotype and gene expression in the lung would be required to provide full validation.

In conclusion, integration of lung-specific eQTL data with GWAS from the same individuals has revealed interesting and potentially causal pathways of airflow obstruction. GWAS have revolutionised our ability to identify gene variants that contribute to susceptibility for common complex genetic diseases but often do not pinpoint the exact genes or mechanisms. Causal pathway analysis involving the joint examination of genetic and genomic data is a vital next step in discovering novel biomarkers and therapeutic targets in airflow obstruction as evidenced in the present study.

### Author affiliations

<sup>1</sup>Institut universitaire de cardiologie et de pneumologie de Québec, Québec, Canada <sup>2</sup>Department of Pathology and Medical Biology, University of Groningen, University Medical Center Groningen, GRIAC Research Institute, Groningen, The Netherlands <sup>3</sup>Department of Genetics and Genomics Sciences, Mount Sinai School of Medicine, New York, New York, USA

<sup>4</sup>Department of Molecular Medicine, Laval University, Québec, Canada

<sup>5</sup>Division of Computational Biomedicine, Bioinformatics Program, Boston University, Boston, Massachusetts, USA

<sup>6</sup>University of British Columbia Center for Heart Lung Innovation and Institute for Heart and Lung Health, St Paul's Hospital, Vancouver, British Columbia, Canada <sup>7</sup>Department of Pathology and Laboratory Medicine, University of British Columbia, Vancouver, British Columbia, Canada

<sup>8</sup>Department of Pulmonology, University of Groningen, University Medical Center Groningen, GRIAC Research Institute, Groningen, The Netherlands <sup>9</sup>Respiratory Division, Department of Medicine, University of British Columbia,

Vancouver, British Columbia, Canada

<sup>10</sup>British Columbia Cancer Agency, Vancouver, British Columbia, Canada
<sup>11</sup>Merck & Co Inc, Rahway, New Jersey, USA

Acknowledgements The authors would like to thank Christine Racine and Sabrina Biardel at the IUCPQ site of the Respiratory Health Network Biobank of the FRQS for their valuable assistance. They also acknowledge the staff at Calcul Québec for IT support with the high-performance computer clusters. At UBC the authors thank the biobank coordinator Dr Mark Elliott. At the Groningen UMCG site Marnix Jonker is thanked for his support in selecting, handling and sending of lung tissues. The authors thank Dr Joshua Millstein of University of Southern California for helpful advice on causality inference methodology.

**Contributors** WT, KH, YB, MLav, PDP, DDS and DSP were involved in the conception and design of the study. MLam, WT, KH, YB, MLav, KSt, JDC, CC, MC, KSh, JCH, C-AB, SL, MEL, ASp, PDP, DN and DSP were involved in the acquisition and/or the analysis and interpretation of the data. MLam, WT, KH, YB, MB, ASa, PDP, DDS and DSP contributed to the writing or the revision of the manuscript.

**Funding** This study was funded by Merck Research Laboratories, the Chaire de pneumologie de la Fondation JD Bégin de l'Université Laval, the Fondation de l'Institut universitaire de cardiologie et de pneumologie de Québec, the Respiratory Health Network of the FRQS, the Cancer Research Society and Read for the Cure, and the Canadian Institutes of Health Research (MOP-123369). YB was a research

scholar from the Heart and Stroke Foundation of Canada and he is now recipient of a Junior 2 Research Scholar award from the Fonds de recherche Québec—Santé (FRQS). DDS is a Tier 1 Canada Research Chair for COPD. ML is the recipient of a doctoral studentship from the Fonds de recherche Québec - Santé (FRQS).

**Competing interests** CAB received a grant from Rosetta Merck. DN is a full-time employee of Merck. DSP received consultancy fees from AstraZeneca, Boehringer Ingelheim, Chiesi, GlaxoSmithKline, Takeda, TEVA and a grant from Chiesi. DDS has served on advisory boards of Almirall, Nycomed, Talecris, AstraZeneca, Merck Frosst, Novartis and GlaxoSmithKline, received grants from AstraZeneca, GlaxoSmithKline and Wyeth, and received honoraria for speaking engagements from Takeda, AstraZeneca, GlaxoSmithKline and Boehringer Ingelheim. JDC received consultancy fees from Metera Biosciences and Immuneering. MB received grants from TEVA, AstraZaneca and Chiesi. WT received a grant from Merck Sharp Dohme, received consultancy fees from Pfizer, and received lecture fees from GlaxoSmithKline, Chiesi and Roche Diagnosis.

Patient consent Obtained.

**Ethics approval** Institutional Review Board guidelines at the three sites.

Provenance and peer review Not commissioned; externally peer reviewed.

### REFERENCES

- Manolio TA. Genomewide association studies and assessment of the risk of disease. N Engl J Med 2010;363:166–76.
- 2 Cho MH, Boutaoui N, Klanderman BJ, *et al*. Variants in FAM13A are associated with chronic obstructive pulmonary disease. *Nat Genet* 2010;42:200–2.
- 3 Cho MH, Castaldi PJ, Wan ES, et al. A genome-wide association study of COPD identifies a susceptibility locus on chromosome 19q13. Hum Mol Genet 2012;21:947–57.
- 4 Pillai SG, Ge D, Zhu G, et al. A genome-wide association study in chronic obstructive pulmonary disease (COPD): identification of two major susceptibility loci. PLoS Genet 2009;5:e1000421.
- 5 Bosse Y. Updates on the COPD gene list. *Int J Chron Obstruct Pulmon Dis* 2012;7:607–31.
- 6 Altshuler D, Daly MJ, Lander ES. Genetic mapping in human disease. Science 2008;322:881–8.
- 7 Thun GA, Imboden M, Ferrarotti I, et al. Causal and synthetic associations of variants in the SERPINA gene cluster with alpha1-antitrypsin serum levels. PLoS Genet 2013;9:e1003585.
- 8 Nicolae DL, Gamazon E, Zhang W, et al. Trait-associated SNPs are more likely to be eQTLs: annotation to enhance discovery from GWAS. PLoS Genet 2010;6:e1000888.
- 9 Hao K, Bosse Y, Nickle DC, et al. Lung eQTLs to help reveal the molecular underpinnings of asthma. PLoS Genet 2012;8:e1003029.
- 10 Lamontagne M, Couture C, Postma DS, et al. Refining susceptibility loci of chronic obstructive pulmonary disease with lung eQTLs. PLoS ONE 2013;8:e70220.
- 11 Sieberts SK, Schadt EE. Moving toward a system genetics view of disease. Mamm Genome 2007;18:389–401.
- 12 Schadt EE, Lamb J, Yang X, et al. An integrative genomics approach to infer causal associations between gene expression and disease. Nat Genet 2005;37:710–7.
- 13 Steiling K, van den Berge M, Hijazi K, et al. A dynamic bronchial airway gene expression signature of chronic obstructive pulmonary disease and lung function impairment. Am J Respir Crit Care Med 2013;187:933–42.
- 14 Campbell JD, McDonough JE, Zeskind JE, et al. A gene expression signature of emphysema-related lung destruction and its reversal by the tripeptide GHK. Genome Med 2012;4:67.
- 15 Butler MW, Fukui T, Salit J, et al. Modulation of cystatin A expression in human airway epithelium related to genotype, smoking, COPD, and lung cancer. Cancer Res 2011;71:2572–81.
- 16 Pavlova A, Bjork I. Grafting of features of cystatins C or B into the N-terminal region or second binding loop of cystatin A (stefin A) substantially enhances inhibition of cysteine proteinases. *Biochemistry* 2003;42:11326–33.
- 17 Rokadia HK, Agarwal S. Serum cystatin C and emphysema: results from the National Health and Nutrition Examination Survey (NHANES). *Lung* 2012;190:283–90.
- 18 Takeyabu K, Betsuyaku T, Nishimura M, et al. Cysteine proteinases and cystatin C in bronchoalveolar lavage fluid from subjects with subclinical emphysema. Eur Respir J 1998;12:1033–9.
- 19 Ishii T, Abboud RT, Wallace AM, et al. Alveolar macrophage proteinase/antiproteinase expression and lung function/emphysema. Eur Respir J 2014;43:82–91.
- 20 Li J, Dai A, Hu R, et al. Positive correlation between PPARgamma/PGC-1alpha and gamma-GCS in lungs of rats and patients with chronic obstructive pulmonary disease. Acta Biochim Biophys Sin (Shanghai) 2010;42:603–14.
- 21 Trian T, Benard G, Begueret H, et al. Bronchial smooth muscle remodeling involves calcium-dependent enhanced mitochondrial biogenesis in asthma. J Exp Med 2007;204:3173–81.
- 22 Remels AH, Gosker HR, Schrauwen P, et al. TNF-alpha impairs regulation of muscle oxidative phenotype: implications for cachexia? FASEB J 2010;24:5052–62.

## Chronic obstructive pulmonary disease

- 23 Wilk JB, Walter RE, Laramie JM, *et al.* Framingham Heart Study genome-wide association: results for pulmonary function measures. *BMC Med Genet* 2007;8(Suppl 1):S8.
- 24 Yanbaeva DG, Wouters EF, Dentener MA, et al. Association of glutathione-Stransferase omega haplotypes with susceptibility to chronic obstructive pulmonary disease. Free Radic Res 2009;43:738–43.
- 25 de Andrade M, Li Y, Marks RS, et al. Genetic variants associated with the risk of chronic obstructive pulmonary disease with and without lung cancer. Cancer Prev Res (Phila) 2012;5:365–73.
- 26 Park JW, Ryter SW, Kyung SY, et al. The phosphodiesterase 4 inhibitor rolipram protects against cigarette smoke extract-induced apoptosis in human lung fibroblasts. Eur J Pharmacol 2013;706:76–83.
- 27 Yang Q, Underwood MJ, Hsin MK, et al. Dysfunction of pulmonary vascular endothelium in chronic obstructive pulmonary disease: basic considerations for future drug development. Curr Drug Metab 2008;9:661–7.
- 28 Morissette MC, Vachon-Beaudoin G, Parent J, et al. Increased p53 level, Bax/Bcl-x (L) ratio, and TRAIL receptor expression in human emphysema. Am J Respir Crit Care Med 2008;178:240–7.
- 29 Bodmer JL, Holler N, Reynard S, *et al.* TRAIL receptor-2 signals apoptosis through FADD and caspase-8. *Nat Cell Biol* 2000;2:241–3.
- 30 He JQ, Connett JE, Anthonisen NR, et al. Glutathione S-transferase variants and their interaction with smoking on lung function. Am J Respir Crit Care Med 2004;170:388–94.