

# Surrogate Endpoints and Drug Regulation: What Is Needed to Clarify the Evidence

*Spencer Phillips Hey, William B. Feldman, Emily H. Jung, Elvira D'Andrea, and Aaron S. Kesselheim*

In the development of new prescription drugs, one company typically owns the drug by virtue of controlling its intellectual property rights. That company also often funds the clinical trials needed to test the drug's efficacy. In many cases, the effect observed in these clinical trials is a change with respect to a biomarker. A biomarker is a characteristic of a disease or a response to an intervention that can be measured via a laboratory test, radiology study, physical exam finding, or other clinical test. Identifying new biomarkers for use in concert with investigational drugs can be a central part of drug development. Yet, unlike drug development, when researchers identify or develop new biomarkers, there is often no single company or stakeholder that owns the biomarker.<sup>1</sup> Consequently, there is no stakeholder in the primary role of generating the evidence to show how a biomarker is useful. Indeed, if the biomarker could be useful across multiple drugs, there may be a collective action problem in that any investment by one stakeholder in validating the biomarker will benefit all the others and consequently, drug developers may not want to invest their limited scientific resources in ways that will help their competitors. There is also no central regulator tracking the evidence for or against a particular biomarker or helping ensuring its proper use.

One key function of biomarkers is to serve as surrogate measures in the pivotal clinical trials testing drugs for regulatory approval. Investigators rely upon surrogate measures because changes to these measures

---

**Spencer Phillips Hey, Ph.D.**, is with the Program On Regulation, Therapeutics, And Law (PORTAL), Division of Pharmacoepidemiology and Pharmacoeconomics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, Massachusetts and Center for Bioethics, Harvard Medical School. **William B. Feldman, M.D., D.Phil.**, is with the Program On Regulation, Therapeutics, And Law (PORTAL), Division of Pharmacoepidemiology and Pharmacoeconomics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, Massachusetts and the Division of Pulmonary and Critical Care Medicine, Department of Medicine, Brigham and Women's Hospital. **Emily H. Jung, A.B.**, is with the Program On Regulation, Therapeutics, And Law (PORTAL), Division of Pharmacoepidemiology and Pharmacoeconomics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, Massachusetts. **Elvira D'Andrea, M.D., M.P.H.**, is with the Program On Regulation, Therapeutics, And Law (PORTAL), Division of Pharmacoepidemiology and Pharmacoeconomics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, Massachusetts. **Aaron S. Kesselheim, M.D., J.D., M.P.H.**, Program On Regulation, Therapeutics, And Law (PORTAL), Division of Pharmacoepidemiology and Pharmacoeconomics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, Massachusetts and the Center for Bioethics, Harvard Medical School.

Table 1

**Extracted from the FDA's Tables of Surrogate Endpoints Used for Drug Approval**

Disease or Use	Patient Population	Surrogate Endpoint	Type of Approval	Drug Mechanism
Asthma	Patients with asthma	Forced expiratory volume in 1 second (FEV1)	Traditional	Corticosteroid; beta-2 adrenergic agonist
Cancer: hematologic malignancies	Patients with B-cell precursor acute lymphoblastic leukemia in first or second complete remission	Minimal residual disease response rate	Accelerated	Mechanism agnostic
Hypercholesterolemia	Patients with heterozygous familial and nonfamilial hypercholesterolemia	Serum LDL-C	Traditional	Lipid-lowering
Type 2 diabetes mellitus	Patients with type 2 diabetes mellitus	Serum hemoglobin A1C	Traditional	Glucose-lowering

Source: US Food and Drug Administration, "Table of Surrogate Endpoints That Were the Basis of Drug Approval or Licensure," available at <<https://www.fda.gov/Drugs/DevelopmentApprovalProcess/DevelopmentResources/ucm613636.htm>> (last visited July 17, 2019).

often occur much sooner than changes to the direct measures of how patients feel, function, or survive. A biomarker that serves as a good surrogate measure can be relied upon to predict an actual clinical outcome months or years before that outcome occurs. For example, when patients are treated with a cholesterol-lowering "statin" medication, like atorvastatin (Lipitor), changes to LDL cholesterol that occur in days can predict changes in the rates of cardiovascular events that may occur years later.

However, not all biomarkers are valid predictors of clinical outcomes when used as surrogate measures in clinical trials. There have been dozens of cases in which drug treatment leading to biologically plausible biomarker changes do not result in predicted clinical changes (or even lead to unanticipated worse clinical outcomes).<sup>2</sup> Unfortunately, the lack of centralized oversight relating to biomarkers can mean that an invalid surrogate (or a surrogate of uncertain validity) may be used as a primary endpoint in pivotal trials and potentially serve as the basis for regulatory approval.

In July 2018, the U.S. Food and Drug Administration (FDA) published two tables of surrogate endpoints that have been used as the basis for past drug approvals, one for adults and one for pediatric indications.<sup>3</sup> By demonstrating how the FDA views the strength of evidence for particular surrogate endpoints, these tables may fill a gap in biomarker oversight and serve as the basis for improving the study and appropriate use of surrogate endpoints. As of April 2019, the table of surrogate mea-

asures for adult indications had 101 rows and the pediatric table had 56 rows. Each row specifies the particular use of a surrogate measure in five fields: (1) disease; (2) patient population; (3) surrogate; (4) type of approval (i.e., traditional or accelerated); and (5) mechanism of action. **Table 1** presents an excerpt from the adult table, showing the rows for asthma, hematologic malignancies, hypercholesterolemia, and type 2 diabetes.

The FDA's plan is to update this list at regular intervals to facilitate discussion with drug developers about appropriate endpoints in pivotal trials.<sup>4</sup> The FDA's website also states that decisions about the acceptability of surrogate measures will be made on a case-by-case basis — taking into account "the disease, studied patient population, therapeutic mechanism of action, and availability of current treatments" — and that surrogate measures are not included in, or will be removed from, the table in accordance with the evolving state of scientific understanding.<sup>5</sup> In other words, the FDA intends this table to serve as a regulatory heuristic — a guide to the FDA's policies about appropriate surrogate measures to use as primary trial endpoints for new drug approvals, but not a replacement for discussion with the FDA about pivotal trial endpoints. In this essay, we will evaluate the extent to which these tables successfully meet this goal and how they can be augmented to improve the drug development process.

## Details of the FDA's Surrogate Endpoint Tables

The tables in their current forms provide some information about biomarkers and other surrogate measures in a systematic, centralized fashion for the first time. The tables helpfully indicate that the context of use for surrogate endpoints does not necessarily extend to an entire disease or condition, but may instead be limited to certain mechanisms or classes of drugs. For researchers interested in developing biomarkers as new surrogate endpoints, this information could serve as a valuable reference point for the kinds of biomarkers and surrogates that the FDA considers to be valid. The tables also usefully highlight the distinction between the kinds of surrogate endpoints that the FDA considers to be valid for full approvals versus accelerated approvals — since the latter require the new drug's manufacturer to conduct additional follow-up trials to prove that their product does indeed have a favorable benefit/risk profile on a clinical endpoint.

Unfortunately, the tables do not specify the clinical or patient-centered outcome for which the biomarker is supposed to be a trial-level surrogate. A valid trial-level surrogate is a measure for which the difference between trial treatment arms with respect to a surrogate outcome has been shown to reliably predict a difference between treatment arms with respect to a patient-centered outcome.<sup>6</sup> To claim that a biomarker — e.g., forced expiratory volume in 1 second (FEV1) — is a valid trial-level surrogate signals that the biomarker is a good predictor of one or more clinical endpoints reflecting how a patient feels, functions, or survives (e.g., reduction in acute asthma exacerbations requiring systemic steroids).

For some diseases, such as cancer, the outcomes for which the biomarker is acting as a surrogate may be obvious (e.g., survival). But in other cases, the real clinical outcomes may be unclear. For example, serum low-density lipoprotein cholesterol (LDL-C) is included as a surrogate endpoint for hypercholesterolemia. But lowering a patient's serum LDL-C is not a surrogate measure (strictly speaking) for hypercholesterolemia. The disease is itself defined in terms of high cholesterol, and therefore lower cholesterol is a direct measure of improvement in the disease. Patients are diagnosed and treated for hypercholesterolemia because high cholesterol is a risk factor for cardiovascular disease, so this is likely the outcome implied by the FDA. But to avoid such assumptions, the patient-centered outcomes should be explicit throughout the tables.

The tables in their current form also lack references to the supporting evidence. Providing such evidence

would extend the table's usefulness by elucidating the quality of evidence that the FDA takes to be sufficient to show that a particular biomarker is a good trial-level surrogate. Showing the evidence is also important since surrogates can link to clinical outcomes in complex ways. For example, hemoglobin A1c is a frequently-used biomarker for glycemic control in the management of diabetes. But for the drug rosiglitazone, evidence suggests that the drug reduces hemoglobin A1c but may paradoxically increase the risk for adverse cardiovascular outcomes. By contrast, two other medications used to treat diabetes — empagliflozin and liraglutide — appear to have positive cardiovascular benefits that are not predicted or explained by their capacity to reduce hemoglobin A1c.<sup>7</sup> Since hemoglobin A1c is not always a reliable trial-level surrogate, other factors need to be considered to determine under what circumstances hemoglobin A1c should be used as a surrogate endpoint in pivotal trials. This further underscores the need to make the supporting evidence base for each row explicit. Safe and reliable use of biomarkers requires that we understand the boundaries of utility around the biomarker's use. Providing links to the evidence would allow the FDA to highlight exceptions to rows in the table that could then be documented and explained.

## Developing Enhanced Surrogate Measure Tables: The Case of FEV1

To address these limitations, we propose that the FDA's tables should include at least one more column to specify the patient-centered outcomes for which the biomarker is believed to be a valid trial-level surrogate. The tables should also provide links to existing meta-analyses (either in the published literature or conducted by the FDA) that provide evidence supporting the validity of the trial-level surrogate. That is, the tables should provide a more comprehensive context-of-use for the surrogate and show that the use is grounded in the existing evidence.

To demonstrate such an expansion of the table, we piloted the process of adding this information for the asthma/FEV1 row. We reviewed the Cochrane Database of Systematic Reviews looking for meta-analyses that would support the implication of the FDA's tables that FEV1 is a valid trial-level surrogate that can be used to support the regulatory approval of corticosteroids or beta-2 agonists for the treatment of adults with asthma.

In January 2019, we searched the Cochrane Library for all reviews containing the terms “asthma” and (“steroid” or “beta-2”) in the title, abstract, or as a keyword; and the term “forced expiratory volume” in any text. This search returned 80 results, which were then

independently screened for relevance (by SPH and ED) and extracted (by SPH) for the PICOS elements (population, intervention, comparator, outcomes, study-types).

We found 19 reviews relevant to the question of whether FEV1 is an appropriate trial-level surrogate for patient-centered outcomes in adult asthma trials testing corticosteroids or beta-2 agonists. The characteristics of this sample are described in **Table 2**. The majority of reviews included in our sample examined treatments for chronic asthma (68%) and included some pediatric populations (78%). All reviews focused on randomized controlled trials (RCTs) as the study design of interest, although 2 (8%) reviews also included “quasi-RCTs” (e.g., comparative trials that used non-random methods of allocation). Corticosteroids and beta-2 agonists were interventions of interest in 15 (78%) and 10 (52%) reviews, respectively. Five reviews (26%) compared different formulations/deliv-

eries of corticosteroids (e.g., oral vs. inhaled). Corticosteroids were the most frequent drug class included as a comparator (63%). Six (31%) reviews involved a placebo comparator.

There were also four comparator drug classes that appeared in only 1 or 2 reviews each: sodium cromoglicate, xanthines, anti-leukotrienes, and muscarinic antagonists. This heterogeneity across comparator classes highlights another dimension of uncertainty with respect to trial-level surrogates. Can investigators assume that a trial-level surrogate’s validity will be constant across all comparison classes? If not, then the comparison classes or mechanisms for which a surrogate may be valid or invalid are another component of a surrogate’s context-of-use that should be made explicit in the table.

In **Table 3**, we list the primary outcomes and the different FEV1 outcome(s) that were measured across our sample of reviews. For primary outcomes, exacerbations requiring systemic steroids was the most common patient-centered outcome (52%), followed by adverse events (42%), hospital admission (26%), quality of life (21%), acute asthma relapse requiring unplanned medical care (21%), symptom resolution (15%), and length of hospital stay (5%). While not an exhaustive list, we believe this set of 7 patient-centered outcomes could provide a plausible place to begin filling out the content for an improved table of surrogate measures by specifying the patient-centered outcomes for which the biomarker may be a valid trial-level surrogate.

Alternatively, the FDA could draw on (or refer the table user to) lists of clinically meaningful, patient-centered outcomes published by professional societies. For example, the American Thoracic Society (ATS) has published a list of endpoints that they recommend for clinical trials of asthma therapy. The ATS’s list of recommended primary endpoints for pivotal trials is similar to our compiled list of clinical outcomes, including symptom-free days, reliever use, exacerbations requiring systemic steroids, and quality of life.<sup>8</sup> Whatever the methodology, it is important that the FDA generate a list of asthma-specific clinical endpoints that reflect consensus among researchers and clinicians about the core outcomes that matter for patients. A surrogate measure that predicts outcomes of minor (or no) importance to patients may not be suitable for use by the FDA when approving new drugs.

Across our sample of Cochrane reviews, we found 4 different outcome measurements that were used to assess an intervention’s effect on FEV1. The most frequent measurements were change in FEV1 compared to baseline (84%), change in percentage of

Table 2

### Survey of Cochrane Reviews Involving FEV1 and Asthma

Characteristics	N	%
<b>Conditions</b>		
Asthma, Chronic	13	68
Asthma, Acute	6	31
<b>Populations</b>		
Adults	19	100
Children	15	78
Post-Emergency Room	4	21
Steroid-naive	1	5
Mild	1	5
<b>Study Types</b>		
Randomized controlled trials	19	100
Quasi-RCT	2	10
<b>Interventions</b>		
Corticosteroid	15	78
Beta-2 agonist	10	52
<b>Comparators</b>		
Corticosteroid	12	63
Placebo	6	31
Beta-2 agonist	2	10
Xanthines	2	10
Sodium cromoglycate	1	5
Aminophylline	1	5
Anti-leukotrienes	1	5
Muscarinic antagonists	1	5



FEV1 predicted (57%), area under the curve (AUC) for change in FEV1 (10%), and change in FEV1 post-exercise (5%). While all of these endpoints are based on common, standardized measurement techniques (using spirometry), this heterogeneity in methods of aggregating population-level effects on FEV1 provides another important lesson about the FDA's tables. A scientifically rigorous claim about a trial-level surrogate requires fully-specified outcomes — i.e., stating the measurement (e.g., forced expiratory volume in 1 second), the technology used to make the measurement (e.g., spirometer), the measurement technique employed with the technology (e.g. how respiratory therapists conduct the test), the metric (e.g., change in liters from baseline), a method of aggregation (e.g., mean), and a time-point (e.g., 3 months).<sup>9</sup> Clearly defining the appropriate context of use for a surrogate endpoint biomarker requires that these six components are specified for the surrogate measure and the patient-centered outcomes.

In **Table 4**, we summarize the results of the 19 Cochrane reviews in our sample using an outcome matrix, arranging the FEV1 endpoints along the x-axis and the clinical endpoints along the y-axis. Each dot in this matrix corresponds to the qualitative concordance between an FEV1 endpoint and a clinical endpoint. For example, a review that included meta-analyses for treatment differences on exacerbations requiring systemic steroids, symptom resolution, and

Table 3

**Primary and FEV1 Outcomes in Cochrane Reviews Involving FEV1 and Asthma**

Characteristic	N	%
<b>Primary Outcomes</b>		
Exacerbations requiring systemic steroids	10	52
Adverse events	8	42
Hospital admission	5	26
Quality of life	4	21
Acute asthma relapse requiring unplanned medical care	4	21
Change in PEF	4	21
Change in percentage FEV1 predicted	3	15
Change in forced vital capacity	3	15
Symptom resolution	3	15
Change in FEV1	3	15
AUC FEV1 change	2	10
Length of hospital stay	1	5
<b>FEV1 Outcomes</b>		
Change from baseline	16	84
Change in percentage predicted	11	57
AUC change	2	10
Change post exercise	1	5

PEF: peak expiratory flow; FEV1: forced expiratory volume in 1 second; AUC: Area under the curve

Table 4

**Outcome Matrix for Survey of Cochrane Reviews Studying FEV1 and Clinical Endpoints in Asthma**

		Clinical Endpoints						
		Exacerbations requiring systemic steroids	Adverse events	Quality of life	Hospital admission	Symptom resolution	Acute asthma relapse	Length of hospital stay
FEV1 Endpoints	Change from baseline	●●●●● ●●○OX	●○○○□ XX	○□X	●○XX	●X	XX	○
	Change after exercise	□						
	Change in % predicted	●○○X	○○X	○X	●X	○X	XXX	
	AUC change	●	○□	○		●		

AUC=Area under the curve

Each marker in this matrix corresponds to results from one Cochrane review. Markers represent the correlation in the direction of effect between the FEV1 endpoint and the clinical endpoint.

- = difference in effects on FEV1 and clinical endpoint significantly in favor of the same intervention
- = significant difference in effect observed for FEV1 or clinical endpoint (but not both); or no significant difference on either outcome
- = significant difference in effects on FEV1 and clinical endpoint, but one favors the experimental intervention and one favors the control
- X = insufficient data for analysis

**Progress in understanding biomarkers and improving their use is not possible unless policymakers recognize the scope of the challenge. Modifying all of the rows of the FDA’s table with the missing content and context is a daunting task, but it will allow the tables to optimally help clinical investigators make evidence-based use of biomarkers in the future.**

change in FEV1 from baseline would produce 2 dots in this figure — one for the concordance of effects on exacerbations and change in FEV1 from baseline; one for the concordance of effects on symptom resolution and change in FEV1 from baseline. We coded the dots in this figure according to four categories of concordance: (1) black circles represent analyses that found a significant difference in effects on the FEV1 and clinical endpoint in favor of the same intervention; (2) white circles represent analyses in which a significant difference in effect was observed for the FEV1 or clinical endpoint but not both, or where no significant difference was observed on either outcome; (3) white squares represent analyses in which significant differences were observed for both the FEV1 and clinical endpoint, but one outcome favored the experimental intervention and one outcome favored the control; and (4) x’s represent cases in which the Cochrane

review reported insufficient data for either the FEV1 or clinical endpoint, and thus it was not possible to examine concordance.

This classification scheme reveals that effects on change in FEV1 from baseline appear to generally concord with effects on exacerbations requiring systemic steroids. Only one review examined the change in FEV1 after exercise and found an opposite direction of effect on the clinical endpoint. For the other two FEV1 endpoints — i.e., change in the percentage of FEV1 predicted; area under the curve (AUC) for change in FEV1 — concordance with clinical outcomes has been more variable. Some reviews found that effects on these endpoints concurred with clinical endpoints, but reviews also found evidence that would seem to disconfirm the utility of these endpoints as trial-level surrogates.

While this kind of outcome matrix is not sufficient on its own to draw definitive conclusions about the validity of a trial-level surrogate — and should not replace a rigorous surrogacy analysis — it does still provide some important insights for how we might improve the FDA’s table. For example, in **Table 5**, we present our modified row of the FDA’s adult table for asthma and FEV1. In the column for surrogate endpoint, we have now added the three specific measures of the surrogate for which the overall weight of evidence across Cochrane reviews in our sample favored its use as a trial-level surrogate. We also added a column for “Clinical Endpoint(s),” which captures all of the clinical endpoints identified by our search, and a column for “Clinical Endpoint(s) Predicted by Sur-

Table 5

**Revised Row of FDA Table of Surrogate Endpoints for FEV1 in Adults with Asthma**

Disease or Use	Patient Population	Surrogate Endpoint	Type of Approval	Drug Mechanism(s)	Clinical Endpoint(s)	Clinical Endpoint(s) Predicted by Surrogate Endpoint
Asthma	Adult patients with asthma	Forced expiratory volume in 1 second (FEV1) • Change from baseline • Change in % predicted • AUC change	Traditional	Corticosteroid; beta-2 adrenergic agonist	• Exacerbations requiring systemic steroids • Adverse events • Quality of life • Hospital admission • Symptom resolution • Acute asthma relapse • Length of hospital stay	• Exacerbations requiring systemic steroids • Hospital admission • Symptom resolution

rogate Endpoints,” which captures the subset of clinical endpoints for which data exist demonstrating an association with the surrogate measure. Though there are other ways of establishing such an association, we have preliminarily filled out this final column with a list of clinical endpoints for which the majority of reviews (with sufficient data) found that effects on one of the FEV1 endpoints concurred with effects on the clinical endpoint. If we imagine this modified table as an interactive, online tool, then the user should be able to “click through” on any row in the table to see the outcome matrix and click on each marker in the matrix to access the corresponding Cochrane review.

### Conclusion

Creation of the FDA’s tables is an important step for addressing an unmet need in the oversight of surrogate endpoints that use biomarkers. These tables could also be used to systematically track and transparently communicate the scientific evidence surrounding surrogate measures. In this analysis, we have outlined ways in which such information could be presented to provide key information on the evidence support the usefulness of each biomarker as a trial-level surrogate. We have also piloted a method for tracking this state of evidence, which further highlights the scientific complexities of biomarker research and development.

Progress in understanding biomarkers and improving their use is not possible unless policymakers recognize the scope of the challenge. Modifying all of the rows of the FDA’s table with the missing content and context is a daunting task, but it will allow the tables

to optimally help clinical investigators make evidence-based use of biomarkers in the future.

### Note

This article was supported by a grant from Arnold Ventures. Dr. Kesselheim also receives support from the Harvard-MIT Center for Regulatory Science and the Engelberg Foundation.

### References

1. A.S. Kesselheim and J. Karlawish, “Biomarkers Unbound — The Supreme Court’s Ruling on Diagnostic-Test Patents,” *New England Journal of Medicine* 366, no. 25 (2012): 2338-2340.
2. T.R. Fleming and D.L. DeMets, “Surrogate End Points in Clinical Trials: Are We Being Misled?” *Annals of Internal Medicine* 125, no. 7 (1996): 605-613.
3. US Food and Drug Administration, “Table of Surrogate Endpoints That Were the Basis of Drug Approval or Licensure,” available at <<https://www.fda.gov/Drugs/DevelopmentApprovalProcess/DevelopmentResources/ucm613636.htm>> (last visited July 17, 2019).
4. *Id.*
5. *Id.*
6. E.L. Korn, P.S. Albert, and L.M. McShane, “Assessing Surrogates as Trial Endpoints Using Mixed Models,” *Statistics in Medicine* 24, no. 2 (2005): 163-182.
7. K.J. Lipska and H.M. Krumholz, “Is Hemoglobin A1c the Right Outcome for Studies of Diabetes?” *JAMA* 317, no. 10 (2017): 1017-1018.
8. H.K. Reddel, D.R. Taylor, E.D. Bateman, L.P. Boulet, H.A. Boushey, W.W. Busse, T.B. Casale, P. Chanez, P.L. Enright, P.G. Gibson, and J.C. de Jongste, “An official American Thoracic Society/European Respiratory Society Statement: Asthma Control and Exacerbations: Standardizing Endpoints for Clinical Asthma Trials and Clinical Practice,” *American Journal of Respiratory and Critical Care Medicine* 180, no. 1 (2009): 59-99.
9. I.J. Saldanha, K. Dickersin, X. Wang, and T. Li, “Outcomes in Cochrane Systematic Reviews Addressing Four Common Eye Conditions: An Evaluation of Completeness and Comparability,” *PloS One* 9, no. 10 (2014): e109400.