Smith et al. Mortality associated with lithium and valproate treatment of US Veterans Health Administration patients with mental disorders. *Br J Psychiatry* (doi: 10.1192/bjp.bp.113.138685)
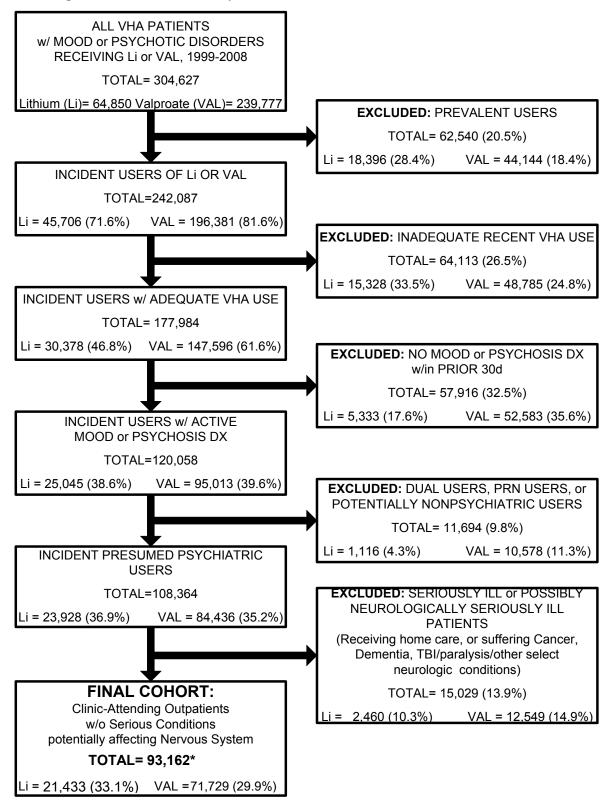
## Online Supplements

**Online Figure DS1  Flowchart of Study Cohort Derivation**

ALL VHA PATIENTS
w/ MOOD or PSYCHOTIC DISORDERS
RECEIVING Li or VAL, 1999-2008

TOTAL= 304,627

Lithium (Li)= 64,850 Valproate (VAL)= 239,777

EXCLUDED: PREVALENT USERS

TOTAL= 62,540 (20.5%)

Li = 18,396 (28.4%)        VAL = 44,144 (18.4%)

INCIDENT USERS OF Li OR VAL

TOTAL=242,087

Li = 45,706 (71.6%)      VAL = 196,381 (81.6%)

EXCLUDED: INADEQUATE RECENT VHA USE

TOTAL= 64,113 (26.5%)

Li = 15,328 (33.5%)        VAL = 48,785 (24.8%)

INCIDENT USERS w/ ADEQUATE VHA USE

TOTAL= 177,984

Li = 30,378 (46.8%)     VAL = 147,596 (61.6%)

EXCLUDED: NO MOOD or PSYCHOSIS DX
w/in PRIOR 30d

TOTAL= 57,916 (32.5%)

Li = 5,333 (17.6%)        VAL = 52,583 (35.6%)

INCIDENT USERS w/ ACTIVE
MOOD or PSYCHOSIS DX

TOTAL=120,058

Li = 25,045 (38.6%)       VAL = 95,013 (39.6%)

EXCLUDED: DUAL USERS, PRN USERS, or
POTENTIALLY NONPSYCHIATRIC USERS

TOTAL= 11,694 (9.8%)

Li = 1,116 (4.3%)          VAL = 10,578 (11.3%)

INCIDENT PRESUMED PSYCHIATRIC
USERS

TOTAL=108,364

Li = 23,928 (36.9%)        VAL = 84,436 (35.2%)

EXCLUDED: SERIOUSLY ILL or POSSIBLY
NEUROLOGICALLY SERIOUSLY ILL
PATIENTS
(Receiving home care, or suffering Cancer,
Dementia, TBI/paralysis/other select
neurologic  conditions)

TOTAL= 15,029 (13.9%)

Li =   2,460 (10.3%)        VAL = 12,549 (14.9%)

**FINAL COHORT:**
Clinic-Attending Outpatients
w/o Serious Conditions
potentially affecting Nervous System
**TOTAL= 93,162***

Li = 21,433 (33.1%)     VAL =71,729 (29.9%)

*excludes 193 patients residing outside the 50 US states

2

**Online supplement DS1  Diagnostic Codes Included in the Cohort**

The databases used in this study were initially developed for use in tracking care delivered to a broad collection of Veterans Health Administration (VHA) patients with depressive or psychotic disorders.  Because of this, a considerable range of diagnostic codes were included during database construction.  To maximize power and because existing literature suggested than any mortality advantages for lithium might span a variety of diagnoses,[13,15,19,20] we decided to retain this broad group of included diagnoses.  Patients could enter the cohort if they had received any one of the following ICD-9 codes in the 30 days prior to lithium or valproate initiation:  (more common) 296.0-296.99, 311, and 295.0-295.9; (less common) 297.0-297.3, 297.8-297.9, 298.0-298.4, 298.8, 300.4, 301.12, 309.0-309.1, and 293.83.   However, virtually all cohort members had received a diagnosis of bipolar I, II, NOS, Depression NOS, major depressive disorder, schizophrenia, schizoaffective disorder, or "other psychoses" within the past 30 days.  Furthermore, these eight diagnostic categories, while still somewhat diverse, were balanced extremely closely in the two matched treatment groups (i.e., within a standardized difference of < 0.019).

**Online supplement DS2  Components of Former User Risk**

"Former user" risks have been proposed as a potential indicator of confounding at baseline[32] or arising during treatment (i.e., confounding arising due to more frequent selection of higher-risk patients in one of the treatment groups to discontinue treatment, also termed here "selection during treatment").[31]  However, in actuality differences in former user risks between treatment groups have the potential to be composed of a complex combination of residual baseline confounding, confounding arising during treatment, differential discontinuation-associated risks, and any difference in persistent effects from treatment.  How often all of these components substantively contribute to former user risk is unclear.

Stated another way, a former user risk (e.g., odds ratio, hazards ratio) of 1.0 is compatible with an absence of confounding, but does not establish this.  Similarity in former user risks cannot establish an absence of confounding if substantial selection during follow-up, discontinuation-associated risks, or persistent effects from active treatment is present.  Since two or more strong biases or effects can potentially co-occur in opposite directions even if the former user risk approximates 1.0, a lack of strong biases can never be definitively concluded on the basis of former user risks alone.  However, former user risks can still have considerable investigative value.  As former user risks gets further from a null value (1.0), it is clear that the presence of confounding, selection, and/or persistence and discontinuation effects becomes increasingly likely, while analyses that achieve former user risks of close to 1.0 may or may not have these substantial effects or biases.

**Online supplement DS3  Aspects of the High-Dimensional Propensity Score Implementation**

Our approach generally followed the overall approach used to develop the original high-dimensional propensity score (hdPS) of Schneeweiss and colleagues,[34] with the following exceptions:

1)  No automated variable construction or selection based on a combined measure of association with exposure and outcome was performed.  Instead initial variables were constructed from most entries in a category (e.g. clinic visits, medications, etc.), except those that were the least common (see below).  Selection was later imposed for the outcome-focused propensity score based solely on associations with outcome, not with exposure.  A few variables were also removed, as described in Online Supplement DS4, relating to very specific measures of past mood stabilizer use which appeared to have the potential to act as instrumental variables.

2)  Although limited screening for covariate prevalence was done (few variables were included if present in <1% of the sample), no limitation was placed on whether that a covariate needed to be present in 5% of the patient sample as in the original hdPS method.[34]  This was because mortality is an infrequent outcome affecting only a small subsample of the cohort, thus even a covariate of low overall prevalence may

contribute to a substantial portion of deaths. Some important covariates judged particularly important *a priori* (e.g. current warfarin prescription, cardiac catheterization in the last 180 days, or age ≥ 80 years old) were included even if present in an overall prevalence of 1% or less.

3) For clearly important variables, a more detailed coding of frequency of occurrence was undertaken than just the absence, presence at < median frequency, presence at greater > median frequency, and presence at > 75[th] percentile frequency categories that were used in the original hdPS method.

4) Greater temporal detail than in the original hdPS method was included for some variables by coding several different time periods for hospitalizations, total provider visits, and other general utilization variables, as well as coding frequency in two time periods for specific clinic visits (0-180 days and 181-365 days) and medications (current prescriptions and recent, but not current, prescriptions [last days' supply ending within the last 180 days]). This strategy was implemented to make information about recent care that might contribute baseline mortality risk less dependent on the exact relationship of the medication initiation date to receipt of medications or services. In addition, this approach might include more detailed information that may be relevant to recent adherence behavior.

**Online Table DS1 Summary of Variables Included in the Initial High-Dimensional Propensity Score (forming the basis of those variables selected for the Outcome-Focused Propensity Score)**

(Prevalence of each balanced within a standardized difference of < 0.019 in the final matched cohort for both the Initial and Outcome-Focused High-Dimensional Propensity Scores)

| Type of Patient Characteristic | Covariates |
|---|---|
| **General Covariates** ||
| Demographics | 10 Covariates (49 indicators) including age (eleven 5-year categories), sex, self-reported race (6), ethnicity, marital status (4), income(6), disability status (4), distance to VHA facility (4), urban/rural hospital location, and fiscal year of medication start (11) |
| Presenting Diagnosis | 9 Covariates denoting psychiatric diagnosis in the past 30 days: Bipolar I, Bipolar II, Bipolar NOS, Major Depression, Depression NOS, Schizophrenia, Schizoaffective Disorder, Other Psychoses, and ≥ 2 of these diagnoses |
| **NonMental Health Covariates** ||
| General Utilization | Total number of current prescriptions, possibly discontinued prescriptions (expired in last 30 days) and recently discontinued prescriptions (expired from 31-180 days), Total number of Provider, Specialist, Surgical, ER Visits, and Inpatient Stays (over from 1- 3 different time periods), number of lab visits in last year, receipt of a flu shot in last year |
| Diagnoses | 41 Variables (44 indicators) relating to diagnoses in the past year, including Total Charlson Comorbidity (CCM) conditions (4 levels), 13 individual CCM Categories (MI, CHF, PVD, CEVD, COPD, Conn Tissue Dz, Peptic Ulcer, Mild Liver Dz, Mod/Sev Liver Dz, DM w/o complications, DM w/ complications, Renal Dz, AIDS/HIV), 11 individual Elixhauser comorbidity categories nonredundant with CCM categories (Arrhythmia ,Weight Loss, Coagulopathy, Pulmonary Circulation Dz, HTN, Valve Dz, Neurodegenerative Dz, Hypothyroid, Obesity, Anemia from Blood Loss, Deficiency Anemia), and 16 additional diagnostic categories (e.g. any fracture, hip fracture, neuropathic pain, back pain, internal injuries, open wounds, etc.), and a Recent Smoking indicator combining Tobacco Dependence diagnosis or treatment) |

| | |
|---|---|
| Current Medications | 54 covariates representing current prescriptions for specific medications or medication categories, including antiarrhythmics, several antibiotics and antihypertensives classes, warfarin, antiplatelet agents, statins, oral diabetes medications, HepC medications, opiate pain medication, low and high dose aspirin, NSAIDS, acetaminophen, inhalers, GI protectants, and other medications |
| Recent Medications | 55 Covariates representing medication/medication categories prescribed in last 180 days but not active on initiation date. Includes 52 of the 54 Current Medication classes, plus vancomycin, antinausea medications, and bandages |
| Hospitalizations | 41 Covariates denoting both the presence of specific types of VHA discharges in the last two years (e.g., Medical ICU, Surgical ICU, Cardiology, Cardiac Stepdown, Telemetry, General Medicine, 8 types of surgical hospitalizations, etc.). An additional set of variables were constructed to indicate the specific type of discharge which constituted the most recent VHA discharge. Additional covariates addressed how recently the latest VHA hospitalization preceded medication initiation and the presence of any discharges against medical advice in the past year |
| Outpatient Providers Visited | 156 Covariates (300 indicators) denoting frequency (typically 0/1/2+ visits) of a large variety of outpatient clinics visited in the last 180 days and in days 181-365 prior to initiation. These include specific medical specialties, specific surgical specialties, anticoagulation clinic, outpatient pharmacy consultation, physical therapy, pacemaker and cardiac catheterization clinics, weight loss clinics, and nonmedical specialty services such work therapy and chaplain visits |
| Diagnostic Tests | 9 covariates (16 indicators) for frequency of tests in prior year, including X-Ray (3), CT/MRI(3), EKG (3), Echocardiogram, Ultrasound, Endoscopy, Nuclear Medicine, PFT, and Angiogram |
| Substance Abuse Diagnoses | 41 variables diagnoses in the past year of alcohol, amphetamine, cannabis, cocaine, opioid, sedative, stimulant, hallucinogen, and other/unspecified substance categories abuse or dependence. For each substance, 4 variables were constructed reflecting the diagnoses categories of abuse, dependence, remission from abuse, and remission from dependence. In addition, variables were constructed to reflect combined drug dependence (with or without opioids), alcohol intoxication, and alcohol psychosis |
| Substance Abuse Treatment | 11 variables (19 covariates), including frequency of individual or group substance abuse treatment in last 180 days or days 181-365 prior to initiation, and current or recent prescription of disulfiram, naltrexone, buprenorphine, or methadone |
| Other Psychiatric Covariates | Numerous covariates including General Mental Health Utilization variables, comorbid psychiatric diagnoses, current psychiatric medications, recent psychiatric medications, number, type and timing of recent psychiatric hospitalizations, recent diagnosed suicide attempts, and types of psychiatric outpatient clinic utilization (e.g., psychiatry, psychotherapy, PTSD-focused, etc.) |
| Aggregate Mortality | 5 indicators denoting age and sex-adjusted state mortality risk, derived from CDC data, grouped into 5 categories (approximate quintiles) |
| VHA Hospital Network Mortality/ Quality of Care | 6 indicators denoting categories of rate of risk-adjusted mortality for the VA Hospital Network (VISN) where patient received care. (Categorization based on data reported in Reference 81) |

**Online Supplement DS4  Derivation of Variables**

This appendix is provided to document for interested readers how the 948 covariates* in the initial high-dimensional propensity score were derived.  These variables form the basis from which the 523 covariates in the outcome-focused propensity score were selected.

### DEMOGRAPHICS AND YEAR OF ENTRY

**Demographics:**  Indicator variables were used for age (< 35 years old, $\geq$ 80 years old, and intervening 5-year age intervals), sex, race, and ethnicity as recorded in Veterans Health Administration (VHA) system.  Where information on race was missing it was imputed using methods previously developed.  Marital status, income, disability status (as indicated by percent of "service connection" of a particular disability), distance to VHA facility, urban/rural location of the facility where patients' obtained care, and fiscal year of medication start was also included.

### UTILIZATION

Utilization variables are derived from VHA clinic stop codes, a set of approximately 500 codes used to categorize each outpatient encounters.  These codes result in classifying care provided into considerably broader categories of care than CPT codes used in other high dimensional propensity scores,[34] reducing the need to consider whether codes should be aggregated or whether information is lost without such aggregation.[77]

**General Mental Health and NonMental Health Utilization:**  We calculated the total number of VHA clinic stop codes relating to encounters with providers (as opposed to Telephone visits, lab tests, etc.) over specific time periods.  We then used multiple indicator variables to categories the frequency with which mental health and nonmental health encounters occurred in the last 7 days before medication initiation and longer time periods over the previous two years.

For general mental health utilization, we also constructed variables reflecting the total number of hospitalizations (as indexed by discharge dates), and variables dividing total MH provider visits into four subtypes (diagnostic interviews, medical management visits, and individual and group psychotherapy visits) over different time periods. For general nonmental health utilization, we also included counts for total nonmental health hospitalizations and the number of surgery clinic and specialist visits (based on stop codes) during particular time periods.  Also, variables were constructed reflecting the total ER/Urgent care visits, lab visits, and presence and absence of a flu shot in the last year (one possible indicator of preventative care).

Finally, for both general mental health and nonmental health utilization, we included indicator variables for the total number of mental health and nonmental health medications.  Different sets of indicator variables accounted for the number of medications that people were receiving on the lithium/valproate start date, the number of medications that they had very recently been taking but for which an active prescription did not exist on the date of lithium/valproate start (termed "Possibly Discontinued"), and the number of medications recently received (within the last 180 days) but not received in the last 30 days ("Recently Discontinued").  A description of the derivation of the covariates designating various specific medications or medication classes which were then summed into these separate counts of total mental health and nonmental health medications is provided in the "Medications" section further below.

*NOTE concerning covariate count:  In the manuscript and here, we use the term "number of covariates" (e.g., 948, 523) to refer to the number of separate, unique quantities balanced through the hdPS-matching. This includes "0 count" quantities for multilevel variables (for these variables, but not dichotomous variables, the number of individuals lacking a positive count for that indicator is a separate quantity, rather than simply information that also can be obtained from the count of individuals scoring "1" for the indicator).  Thus, depending on whether one is considering distinct patient characteristics, number of variable terms entered, or number of unique values balanced

by the model, counts vary.  For example, 546 unique patient qualities were modeled in the initial propensity score using 788 total variables resulting in 948 distinct, unique values, termed here as "covariates".

**Mental Health and NonMental Health Outpatient Utilization:**  Clinic stop codes were classified with indicator variables to reflect whether a patient had attended no visits of that type, a single isolated visit, or repeated visits (2 or more visits of that type) within a time period.  The two time periods examined were the last 180 days prior to lithium/valproate start, and the prior 181 to 365 days before lithium/valproate start.  For mental health outpatient utilization, visits were classified as occurring with psychiatrists, psychotherapists, in the general mental health clinic, primary care behavioral health clinic, substance use disorder clinic, or the Health Care for Homeless Veterans clinic, with additional indicators for visits involving group treatment.

A much greater variety of stop codes exists for nonmental health outpatient utilization.  We chose all stop codes appearing for ≥ 5% of either treatment group in either the last 180 days or days 181 to 365 prior to medication start and other, lower prevalence clinic stop codes thought *a priori* to be of importance (e.g., pacemaker clinic, etc.)

In addition, nonmental health stop codes also were also used to construct the diagnostic testing module described below.

**Mental Health and NonMental Health Hospitalizations:**  The VHA uses approximately 90 bedsection codes to classify hospitalizations by the type of care received (e.g., Specialty of Ward where patient is admitted).  The 30 bedsections that relate to mental health hospitalizations were classified into 4 larger classes: Psychiatric-focused hospitalizations, Substance Abuse-focused hospitalizations, Residential/Day program, and Domiciliary Program (longer-term housing).  With regard to bedsection codes for NonMental Health hospitalizations, a few codes were consolidated when counts were observed to be particularly low (e.g., dermatology bedsection discharges), but in most cases a simple indicator variable was developed to reflect either that the patient's most recent hospitalization had been that bedsection, or that any of their hospitalization bedsections in the two years prior to medication start had been in that bedsection.  These latter variables were constructed both as a measure of overall disease burden (of conditions of a severity requiring hospitalization), because for some progressive conditions earlier hospitalizations or diagnoses can actually reflect worse health prognosis,[78] and because failing health is an obvious risk factor for mortality.  These variables included ICU Bedsections, Step Down Bedsections, Telemetry Bedsections, General Medicine Bedsections, Specialty Medicine (e.g., Neurology, Cardiology) Bedsections, Surgery Bedsections, etc.

Because mortality risks with relation to mental and nonmental health hospitalizations appear to be time-dependent, we coded hospitalizations to whether hospitalizations of any particular type were present in the last 2 years, what the nature of the most recent hospitalization was, and whether the most recent hospitalization was focused on Mental Health or NonMental Health conditions.   We also constructed multiple indicators to reflect the timing of the latest psychiatric discharge date relative to medication initiation to partially reflect the severity of the patients' psychiatric condition, was well as also having an indicator of the total number of nonmental health VHA hospitalizations total in the last year.

## DIAGNOSES

**Comorbid Mental Health and NonMental Health Diagnoses and Indicating Diagnoses:**  Indicator variables were used to reflect a variety of specific mental health diagnoses given in the past year, based on ICD-9.  We required all cohort members to have VHA service use in the past year as well as a prior year, so use of this past year time period helped maximize information about what diagnoses a patient likely actually had.  The one exception were the diagnoses presumed to serve as an indication for lithium/valproate treatment (mood or psychotic diagnoses), for which we required the diagnosis to be entered in the last 30 days.  This briefer period was used in order to maximize the likelihood that this was the reason the patient was receiving lithium or valproate.

Nonmental health diagnoses were aggregated into larger categories based on the comorbid illness categories that make up the Charlson Comorbidity Index and the Elixhauser Comorbidity Index, as per a classification procedure developed for use with administrative databases.[79]  For the Charlson index categories, the following 13 (out of the total 17) comorbidity categories were used:  Myocardial infarction, Congestive Heart

Failure, Peripheral Vascular Disease, Cerebrovascular Disease, Chronic Obstructive Pulmonary Disease, Connective Tissue Disease, Peptic Ulcer, Mild Liver Disease, Moderate or Severe Liver Disease, Diabetes Mellitus without complications, Diabetes Mellitus with complications, Renal Disease, AIDS/HIV Infection.

Elixhauser Comorbidity categories were also included, based on the same reference,[79] when these categories were judged not to overlap with the Charlson index categories. The eleven categories included were: Arrhythmias, Weight Loss, Coagulopathies, Pulmonary Circulation Disease, Valvular Disease, Neurodegenerative Diseases, Hypothyroidism, Obesity, Anemia from Blood Loss, Deficiency Anemia, and Hypertension with or without Complications.

Multiple indicators were also included to reflect the total number of Charlson Comorbidity Index conditions, considering all diagnoses received in the past year.

In addition, indicators for other injury-related and a few specific diagnoses related to progressive neurodegenerative or autoimmune conditions and pain diagnosis were included. Finally, an aggregated smoking indicator was included in this category. Tobacco dependence is recognized as being underdiagnosed in VHA administrative/clinical coding, so we constructed a "recent smoking" variable which assumed a value of "1" if a patient had any of the three in the past year: a diagnosis of Tobacco Dependence, at least one visit to a smoking cessation clinic, or prescription of nicotine replacement therapy or varenicline.

**Comorbid Substance Abuse Diagnoses:** Seven categories of legal/illicit substance use (alcohol, amphetamine, cocaine, marijuana, opioids, sedatives, other substances) were coded as four different indicators reflecting diagnoses received in the last year: dependence on that particular substance, abuse of that particular substance, remission from dependence of that substance and remission from abuse of that substance. The eighth category, hallucinogens, was coded as only 3 indicators (dependence, abuse, and remission from dependence) because there were insufficient numbers of patients ($\leq 5$ in one of the treatment groups) diagnosed with remission from hallucinogen abuse in the past year. In addition, indicators were included for combined substance dependence and remission from combined substance dependence, including separate indicators denoting whether this combined dependence included opioids or not. Two indicators were also included for "unspecified" substance dependence. Lastly, indicators were included in this category for alcohol intoxication (both a narrow and broad definition) and alcohol or drug psychoses.

**Recent Nonfatal Suicidal Behavior Diagnoses:** Episodes of nonfatal suicidal behavior, especially those occurring recently, are a strong risk factor for suicide[80,81] which might plausibly be associated with nonsuicide mortality if any miscoding of suicide deaths occurred.[82] However, there are concerns that outpatient diagnoses sometimes may reflect a history of more remote suicidal behavior rather than reflecting behavior always occurring close to the time the diagnoses were entered. To address these concerns a hierarchy was imposed to avoid double-counting of nonfatal suicide behavior episodes between types of inpatient, or inpatient and outpatient diagnoses. Indicator variables were developed reflecting the occurrence of a diagnosis of an episode of nonfatal suicidal behavior over the last 30 days, days 31 to 180 and days 181 to 365 prior to lithium/valproate start. This approach is expected to result in only an approximate indicator of recent suicide attempts, since a patient could have two separate attempts within a time period that were diagnosed in different settings, and this occurrence would not be reflected in our coding scheme. In addition, the same attempt, if a diagnosis occurred close to the end of a time interval in one setting (e.g., during a non-MH hospitalization), may have been rediagnosed in a second setting in the next time interval. Thus this single behavior episode would appear as two distinct episodes in our coding scheme, not one. Some imprecision of this type is likely unavoidable.

Despite such uncertainties, we felt it was important to incorporate this information when available in our high-dimensional propensity score. Similarly, it was considered important to maintain a distinction concerning the setting of the nonfatal suicidal behavior diagnosis, since an episode diagnosed in a non-mental health hospitalization is likely to be, on average, considerably more serious than diagnoses simply recorded as outpatient diagnoses. It should be recognized that in general diagnoses of nonfatal suicidal behavior are specific but very insensitive,[82] although sensitivity is expected to increase for inpatient diagnoses compared with outpatient (another reason that we made this distinction).

**MEDICATIONS**

**Current and Recent Mental Health Medications:** Mental health medication prescriptions active at the time of lithium/valproate start or recently filled (within the last 180 days) were designated into general classes by 24 indicator variables, using a classification system previously developed. This system already used multiple categories to index antidepressants. For this study we also classified second generation antipsychotics into individual medications (clozapine, olanzapine, risperidone/paliperidone, quetiapine, aripiprazole, ziprasidone). Such an enhanced classification was important given the differential impacts of these medications on mortality and general health (obesity, diabetes, etc.) risk. An identical number of indicator variables were used to reflect recent but not current prescriptions of medications from these same classes, designating receipt of one or more prescription of that type of medication in the last 180 days in the absence of a prescription whose days' supply includes the start date for lithium/valproate treatment.

For nonmental health medications, a system was developed using medication class code information assigned by the VHA by the VHA national formulary. Along with individual medication names and codes, the VHA categorizes every medication administered from the pharmacy into one of more than 1000 classes of medication denoted by a 5 character "medication class" codes, which is further organized into larger 3-character class code categories (using the first 3 digits of the 5 character code). We took advantage of this classification as an approach to logically aggregate prescriptions for related medications (e.g., different loop diuretics were aggregated into the category "Loop Diuretics"). In many cases, we used the broader 3 character code, but in other cases, in which further distinctions were judged important, specific medications were sometimes coded individually (e.g., warfarin). This process condensed the approximately 1000 VHA medication classes used by our cohort down to approximately 225 medication categories. Then all the medication categories present with a prevalence of $\geq 5\%$ in either treatment group (reflecting number of patients with at least one prescription in the last 180 days, or with a current prescription on start date of lithium/valproate) were included. Finally, any medication categories of $< 5\%$ prevalence but $> 1\%$ prevalence that were judged *a priori* particularly relevant to mortality risk (e.g., warfarin, digoxin, etc.) were included as propensity score covariates (this eliminated many of the medication categories, reducing the 225 categories down to 54-55 medication categories [covariates] of clear relevant to our treatment groups).

Indicators for "Current" medication categories required the patient to have an active prescription with a number of days of medication supplied that included the start date of lithium/valproate, while indicators for "Recent" (but not current) revised medication categories required the patient to have had at least one prescription filled in the last 180 days but no active supply at time of lithium/valproate start. A list of some important nonmental health medications or medication categories coded as "Current" and "Recent" medications is included in Online Table DS1. The "Possibly Discontinued" distinction (for medications whose supply ended in the last 30 days prior to lithium/valproate initiation) that was used for the variables reflecting *total* nonmental health and mental health medications (see "Utilization" above) was not retained for these covariates designating individual medication categories. Instead, medications prescribed in this fashion, with their supply ending in the last 30 days prior to lithium valproate initiation, were classified as "Recent" medications, along with the medications with a prescribed number of days of medication supplied that had ended earlier in the180 day period prior to medication initiation).

In the rare cases when fewer than 5 individuals within a treatment group after matching had received medications of a particular category currently or recently, this category was either removed from the propensity score model or consolidated with other medication classes. This resulted in small differences, for instance, in the number of nonmental health current medication (54 variables) versus recent medication categories included as propensity score covariates (55 variables).

**Prior Mood Stabilizer Treatment History:** Although we sought to identify incident users through the requirement of a "clean period", some patients, although a clear minority ($\leq 36\%$ of a treatment group), had had more remote treatment with either mood stabilizers of any type or specifically with lithium or valproate. Two indicator variables were included, reflecting past history of treatment with mood stabilizers in general and past history of treatment with lithium or valproate specifically. Additional variables were considering indexing whether the patient had ever

received the same or opposite mood stabilizer previously or whether the patient's most recent past prescription had been the same or opposite mood stabilizer, but these variables displayed a somewhat strong correlation with which medication (lithium or valproate) the patient initiated in the study. Because these variables also related to past treatment, rather than current treatment, it was judged that these variables, while likely helping to explain assignment to lithium or valproate, would be prime candidates to act as instrumental variables, potentially adding bias[38] to our analysis. These variables were not included in the final high-dimensional propensity score model.

**OTHER**

**NonMental Health Diagnostic Testing:** Clinic stop codes reflecting diagnostic procedures over the last 180 days and days 181 to 365 prior to lithium/valproate start were used to construct indicators of the frequency of diagnostic tests over the past year: X-Rays, CT or MRI scans, EKGs, Ultrasound, Echocardiograms, Endoscopy, PFTs, Nuclear Medicine, and Angiograms (for Angiograms, tests were divided as occurring within the last 180d days and in days 181 to 365 prior to lithium/valproate start).

**Geographic All-Cause Mortality Risk:** Indicator variables were constructed to classify patients into 5 categories (approximate quintiles) of age-adjusted regional (state-level) mortality risk, based on publically available data from the Centers of Disease Control for the years 2000 and 2007.[83] Because these statistics would include the deaths of Veterans occurring in this period, there is the potential for control for "predictors" that include outcome-related information, but this bias is expected to be exceedingly small, given the large number of deaths that occurred across these states over eight years, and the fact that our sample accounted for less than 600 nonsuicide deaths over that period. A geographic all-cause mortality indicator was included to guard against the possibility of regional differences in prescribing patterns creating a spurious association between treatment and mortality.

**VHA Hospital System (VISN) Mortality Risk:** Indicator Variables were used to classify patients into 6 categories of risk-adjusted all-cause mortality risk (based on age, gender, Charlson Comorbidity Index, perceived physical health, and perceived mental health), using information from VA surveys administered in 1998 and 1999.[37] Although this information is most accurate for the very beginning of the study period (mid-1999), it was judged that having some indicator of both mortality variation among VA Hospital Systems and potential quality of care in general would be useful to help limit possible spurious associations between a treatment and mortality due to general tendency for hospitals providing higher or lower quality care to have providers who favored one or the other treatment.

Three additional variables were included to help balance the extensiveness of pharmacy records among our recipients: any prior use of VA pharmacy, use > 180days prior to LI/VAL start, and use > 365d prior to LI/VAL start.

**Selection of Covariates for the Outcome-Focused Propensity Score:** Covariates were selected with an association with mortality of +/- 20 percent (odds ratio of $\geq 1.2$ or $\leq 0.83$), as has been done previously.[39] Determining whether a dichotomous variable has a 20% association with mortality is generally straightforward, however for covariates with more than two possible levels (e.g., age), determining which variables are included or excluded becomes more of a matter of judgment. Either highly restrictive (requiring all categories of the variable to have an association with mortality of $\geq 20\%$) or highly permissive criteria (requiring only 1 category of the multilevel variable to have an association of $\geq 20\%$) could be envisioned. We adopted a compromise approach in which multilevel variables were included in the outcome-focused propensity score only if a majority of level of that variable had a +/- 20% odds ratio association with nonsuicide mortality, except for a very few limited exceptions.

**Online Supplement DS5  Initial Propensity Score Results**

Matching upon the initial propensity score produced results that appear to be consistent with some degree of "amplified confounding".[40,41]  For this reason we chose to report the outcome-focused propensity score results throughout the manuscript as likely more unbiased.  However, for completeness, we report the initial propensity score results here (Online Table DS2) and compare these results with Tables 4 and 5 of the manuscript.

**Online Table DS2 Risk of Nonsuicide Mortality (Intent-to-Treat Cohort, Initial Propensity Score-Matched)**

| Time Period | Hazard Ratio (Lithium/Valproate) | | |
| --- | --- | --- | --- |
| | Intent-to-Treat | During Initial Exposure (As-Treated) | Subsequent Nonexposure (Former User) |
| 0-90 Days | 0.72[a] (0.55-0.95) | 0.81 (0.57-1.14) | 0.67 (0.35-1.25) |
| 0-180 Days | 0.97[b] (0.82-1.15) | 0.83 (0.61-1.12) | 1.19 (0.80-1.77) |
| 0-365 Days | 0.87[c] (0.77-0.97) | 0.77 (0.57-1.03) | 0.84 (0.65-1.07) |

Comparing the former user values between the initial and outcome-focused propensity score-matched analyses suggests that the outcome-focused propensity score is less confounded than the initial propensity score:  at 90 and 365 days the former user hazard ratios are closer to 1.0 for the outcome-focused propensity score than the initial propensity score (central estimate HRs 0-90 day: 0.88 (outcome-focused score) versus 0.67 (initial score); 0-365 days: 1.02 (outcome-focused score) versus 0.84 (initial score).  This pattern does not hold for 0-180 days, but in this case the results are consistent with potential discontinuation risks being attenuated by the presence of greater confounding for the initial propensity-score matched cohort in the direction of better outcomes for lithium (central estimate HR = 1.19 versus 1.54 for the outcome-focused score).  In addition, the initial propensity score intent-to-treat and former user hazard ratios over 0-90 and 0-365 days are more similar to the hazard ratios observed prior to propensity score matching than the outcome-focused propensity score hazard ratios (Online Supplement DS7).  This suggests greater residual confounding for the initial propensity score matched cohorts.  Since this analysis includes more covariates, this suggests a greater amplification of unmeasured/incompletely measured confounding as some have suggested can occur with control of measured covariates not substantially associated with outcome.[38,40] Interestingly, the as-treated initial propensity score results are not closer to the unmatched results than the outcome-focused results, a finding that suggests some contribution from random error or that the effects of amplified confounding/less important covariates may warrant more theoretical or empirical investigation.

Also of note, the former user risk of HR = 0.84 over 0-365 days actually exceeds in magnitude the intent-to-treat risk estimate.  This same pattern is observed over 0-90 days.  This pattern of risk suggests that the initial propensity score intent-to-treat risks may be largely or even entirely explained by confounding, and that the former user risks is made up of substantial confounding combined with an additional element (e.g., random error).

**Online Supplement DS6  Propensity Score Matching Details**

Our propensity score matching was performed using greedy-matching involving freely available SAS code from the Mayo Foundation for Medical Education and Research[36] as well as SAS code from Fairies et al., Chapter 3 of the SAS Press book "Analysis of Observational Health Care Data Using SAS".[43]

Because it is not part of this published code, we did not trim our propensity score cohorts to a "Common Support Area" prior to matching.  Perhaps due to the large preponderance of patients initiating valproate compared to lithium and the wide, overlapping propensity score distribution for both medications, very few patients fell outside the "Common Support Area".  This is reflected by the fact that use of fairly standard 0.2 propensity score logit calipers resulted in a narrow propensity score range while including virtually all lithium-treated patients.  We more precisely established this for our highly similar analysis of suicide mortality.  This analysis also involved nearly complete matching of lithium-treated patients, and these patients were matched using a propensity score that included 98% of the covariates included in the initial propensity score for this study.  For this similar analysis, we established that exceedingly few patients fell outside of a Common Support Area (only 0.05% [lithium-treated patients] to 0.12% [valproate-treated patients] of the entire unmatched cohort).

**Online Supplement DS7  Mortality by Treatment in the Unmatched Cohort and Implications for Confounding**

See online Table DS3 below.  Of note, for each effect estimate the outcome-focused high-dimensional propensity-score matched cohort produced estimates in the direction of reducing the effect sizes from the unmatched analysis each of which indicated a stronger association with worsened outcomes among patients initiating valproate.  (However, this change in estimates between the unmatched and matched analyses was minimal for the as-treated effect estimate).  These findings suggest confounding in the overall cohort is in the direction of patients who are less medically ill preferentially receiving lithium (i.e., patients having less risk of nonsuicide mortality at baseline, prior to treatment initiation).  As a result, the unmatched associations show stronger effect sizes favoring lithium treatment than the matched analysis.

**Online Table DS3 Risk of Nonsuicide Mortality over 365 days, by Treatment**

| | Hazard Ratio (Lithium Versus Valproate) (95% Confidence Interval) | | |
|---|---|---|---|
| | | Stratified by Exposure Status | |
| Analysis (Incident Users) | Intent-to-Treat Sample | During Initial Exposure (As-Treated) | During Subsequent NonExposure (Former Users) |
| Unmatched Cohort | 0.74[a] (0.65-0.84) | 0.58[a] (0.45-0.74) | 0.77[a] (0.63-0.94) |
| Outcome-Focused High-Dimensional Propensity-Score Matched Cohort | 0.92[b] (0.82-1.04) | 0.62[b] (0.45-0.84) | 1.02[b] (0.79-1.32) |

[a] P values are: p < 0.0001 (Intent-to-Treat), p < 0.0001 (Current Users); p = 0.011 (Former Users)
[b] P values are:  p = 0.173 (Intent-to-Treat); p = 0.002 (Current Users); p = 0.888 (Former Users)

**Online Table DS4 Persistence with Treatment: Censoring of Patient Cohorts at 90, 180 and 365 days (outcome-focused Propensity Score-Matched Cohort)**

| Treatment Status | 0-90 Day Follow-up | | | | 0-180 Day Follow-up | | | | 0-365 Day Follow-up | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Lithium | | Valproate | | Lithium | | Valproate | | Lithium | | Valproate | |
| | N | (%) | N | (%) | N | (%) | N | (%) | N | (%) | N | (%) |
| Still Receiving Initial Treatment | 9967 | 46.8 | 9852 | 46.3 | 5012 | 23.5 | 5086 | 23.9 | 1621 | 7.6 | 1723 | 8.1 |
| Discontinued Initial Treatment | 10501 | 48.7 | 11046 | 51.3 | 15220 | 71.5 | 15699 | 73.8 | 18445 | 86.7 | 18974 | 89.1 |
| Initiated Opposite Mood Stabilizer[b] | 773 | 3.6 | 321 | 1.5 | 987 | 4.6 | 415 | 2.0 | 1133 | 5.3 | 470 | 2.2 |
| Suicide Mortality | 15 | 0.07 | 15 | 0.07 | 17 | 0.08 | 16 | 0.08 | 18 | 0.08 | 20 | 0.09 |
| Nonsuicide Mortality[c] | 32 | 0.15 | 54 | 0.25 | 52 | 0.24 | 72 | 0.34 | 71 | 0.33 | 101 | 0.47 |

[a] n = 21288 propensity-score matched pairs

[b] i.e., switched directly from lithium to valproate, or directly from valproate to lithium

[c] Because this Table row reflect patients censored due to nonsuicide death occurring before any treatment impersistence, they are equivalent to the "as-treated" counts

**Online Supplement DS8   Expected versus Observed Decrease over Time of Intent-to-Treat Effect Sizes**

It is important to note that if a genuine medication effect exists during active treatment, it can usually be expected that *intent-to-treat* estimates will progressively  weaken over time if the analysis is largely unconfounded. This occurs because fewer and fewer of the intent-to-treat sample over time remain on the treatment that provides the genuine effects.  (This pattern can be different in studies in which patients predominantly switch, rather than discontinue, treatment).

The pattern observed in this study, however, is somewhat different.  First, while the intent-to-treat central estimate does move towards the null (and become nonsignificant) as followup time increases to 0-180 days (HR = 0.97, [0.82-1.15]), the central estimate moves away from the null over 0-365 days (HR = 0.92, [0.82-1.04]). Second, marked changes in the central estimate of risk among former users is observed, becoming substantially elevated and statistically significant at 0-180 days (HR=1.54, [1.01-2.37]), but then returning almost to an entirely null value over 0-365 days (HR = 1.02, [0.79-1.32]).

The lack of statistical significance for many of these estimates limits interpretation, because random variation can explain many of these findings.  Nevertheless, integrating the pattern of central estimate findings appears to indicate that the most parsimonious interpretation of our study (although far from conclusive) would suggest that our results include both genuine, although transient, increased mortality risk occurring shortly after lithium discontinuation, and residual confounding biasing against valproate (i.e., towards the observation of lower risks with lithium).  The possibility of transient "emergent" risks occurs after lithium discontinuation (compared to

valproate) discontinuation is suggested by the magnitude of the intent-to-treat estimate change from 0-90 days to 0-180 days and by the statistically significant former user mortality risks over 0-180 days. The magnitude of the intent-to-treat estimate change from 0-90 days to 0-180 days, since the 0-180 days include the 0-90 day findings, suggest that the overall intent-to-treat central estimate hazard ratio would actually show increased risk over 91-180 days among all patients who had originally initiated lithium. Interestingly, this is the general pattern of risk that our research team observed in a related study of suicide mortality drawn from the same unmatched cohort. Our secondary analyses risk suggests a possible reason for this increased risk associated with lithium initiation: a sufficient number of former users experience some substantial mortality risk associated with lithium discontinuation to potentially outweigh any benefits experienced by the minority of patients who remain on initial lithium treatment during this period. The risks upon lithium discontinuation over 0-180 days deserve particular attention not only because they are statistically significant, but also because they appear to be in the opposite direction of baseline confounding. Both the risks observed prior to matching (more strongly in favor of lithium treatment [online supplement DS7]) and the tendency for both the intent-to-treat and former user hazard ratios to move in the direction of favoring lithium treatment as follow-up time is extended to 365 days suggests that an underlying confounding bias favor lithium exists in our study.

Again, while none of these conclusions can be made definitively, this interpretation suggests that more weight should be put on the significantly increased risk among former users over 0-180 days, both because of its statistical significance and because the association appears to be observed in contrast to some level of confounding (at least baseline confounding) biasing against observing this association. In addition, caution placed on conclusions about lower risks associated with lithium active treatment and the overall balance of intent-to-treat risks and benefits, given that estimates of risk may be biased by at least some confounding in the direction of favoring lithium treatment. This is particularly true for conclusions about the net balance of risks and benefits from initiating lithium compared to valproate treatment over the first 365 days of treatment. These intent-to-treat associations are both not statistically significant and, as discussed further in the next supplement, would require only relatively small amounts of residual confounding for the presumed unconfounded central estimate hazard ratios to favor valproate, rather than lithium.

**Online Supplement DS9  Possible Implications of Even Relatively Modest Residual Confounding on Assessment of Overall Comparative Benefit of Lithium versus Valproate**

As an example, at 180 days the intent-to-treat HR (central estimate) of 0.97 suggests that residual confounding of only approximately 3-4% lower baseline hazard of mortality at baseline (prior to medication initiation) among patients initiating lithium would be sufficient to yield a central estimate of increased, rather than decreased, mortality risk among all lithium initiators over the first 180 days of treatment. Slightly more confounding would be necessary to yield a central estimate of increased risk over 0-365 days. Statistically, the confidence intervals already preclude definitive conclusions concerning net harms or benefits of lithium compared to valproate for these periods. Notably, the intent-to-treat associations over 0-180 days includes the period of strong, statistical significant intent-to-treat associations between lithium initiation and lower mortality risks from 0-90 days. This implies despite the strong (and statistically significant) protective association between lithium and decreased overall mortality in the first 90 days, increased risks over the 91-180 day period must be sufficient to largely eliminate this association.

Of course, since the increased risks are apparent almost exclusively among patients who have discontinued lithium treatments, our data also suggests that if persistence with lithium could be increased, the balance of benefit to risk for lithium compared to valproate might shift further in the direction of favoring lithium.

**Online Supplement DS10  Evidence for and against Substantial Confounding Arising During Treatment**

On theoretical grounds, confounding arising during treatment (also referred to here as selection during treatment) could plausibly explain the reduced risks in as-treated patients and enhanced risk among patients

14

discontinuing lithium.  This would occur if a greater number of medically ill patients had their lithium treatment stopped because of their deteriorating condition than medically ill patients receiving valproate.  In essence, this differential selection would serve to transfer a greater number of high-risk "as-treated" individuals receiving lithium treatment than receiving valproate treatment to the category of "former users" having discontinued treatment.

However, this phenomenon, unless the selection was based primarily on adverse medical risks caused by the medications themselves, would not easily explain the changes in intent-to-treat risks, nor the consistency of the as-treated risks observed.  Finally, this possibility is also rendered less likely by the very similar rates of medication discontinuation observed between the treatment groups receiving the two medications over time.  However, as others have pointed out, this line of reasoning is not firmly conclusive since patients may discontinue medications at the same rate but for different reasons.[31]

Furthermore, a simple model combining both confounding at baseline and arising during treatment (i.e., positing no medication effects on either risks during active treatment or after discontinuation) does not appear sufficient, since in such a model, in order to explain the intent-to-treat findings, confounding would have to change direction from 0-90 and 91-180 days, and then change direction again over 181-365 days, to explain the 0-90, 0-180, and 0-365 day intent-to-treat estimates observed.

It is possible, of course, that random variation does contribute to the 90, 180 and 365 day estimates and perhaps enhances the differences between them, producing spurious changes in direction of the estimates.  However, the probability random error explains the difference between the 0-90, 0-180, and the 0-365 day intent-to-treat estimates entirely (i.e., in the absence of residual confounding or genuine medication effects during treatment or upon discontinuation), or to the differences between the 0-90, 0-180 and 0-365 day former user risks entirely, would be considerably less than 50%.

Thus, the simplest consistent interpretation of the outcome-focused results is that some level of residual baseline confounding biasing against valproate persists in the 0-365 day analyses, although it is not the only contributor to the risk estimates.  When dissociation-associated risks and/or (less likely) effects from confounding arising during treatment weaken from 0-180 days to 0-365 days, any residual baseline confounding might then serve to "pull" the former user risk much closer to the null over a relatively short period, and contribute an overestimation in the intent-to-treat estimates of the amount of decreased mortality risk associated with lithium initiation.


## Online Supplement DS11  Value of Lithium Treatment Persistence

If, as discussed above, the possibility of substantial confounding cannot be fully excluded, then an important implication results:  it becomes difficult to assess whether patients initiating lithium are at greater or lesser overall mortality risk when the combined impact of both the possibility of possibility of reduced mortality risks during active treatment and increased mortality risks after discontinuation are considered.  This is especially true because the margin of beneficial intent-to-treat association observed at 0-180 and 0-365 days is generally small (and statistically nonsignificant).  That is, as pointed out in online supplement DS8, even relatively small amounts of residual confounding (for the 0-180 day and 0-365 day analyses, respectively) would be sufficient to conceal any overall hazardous treatment effects associated with lithium compared to valproate.  For this reason, we note in the manuscript that further research is clearly needed and caution should be exercised regarding any judgments of whether greater or lesser lithium use would be desirable.

However, despite this uncertainty, at least one conclusion can be made with greater confidence: unless sufficient confounding exists to conceal an actually hazardous relationship between lithium and mortality during *active treatment*, then a clear benefit would exist for maximizing persistence with lithium treatment once initiated.  Our data, being nonrandomized, cannot exclude this possibility, although such founding would have to be quite substantial (i.e., considerably more substantial than the confounding sufficient to conceal a net overall, intent-to-treat hazard for lithium)..  However, if lithium is associated with reduced mortality risks or at worse a neutral effect on mortality (relative to valproate) during active treatment, then our data would support benefits to increasing persistence with lithium treatment, once initiated.  Put another way, regardless of whether the *overall* impact of lithium over the entire follow-up period leads to lesser or greater mortality risks than initiation of valproate, as long

as active treatment with lithium is not associated with an increased, rather than decreased, mortality risk, then our data supports the benefits of improved persistence with lithium once initiated. Emphasizing treatment persistence would have the effect of increasing any benefits experienced during active treatment and reducing any risks resulting from discontinuation. Regardless of which risk predominates, once a decision is made to initiate lithium, efforts to boost treatment persistence when feasible appear likely to benefit the patient.

### Online Supplement DS12  Challenges in Completely Modeling Important Risk Factors

It is important to recognize there is inherent difficulty in capturing a fully desirable amount of information concerning some types of variables. Hospitalizations prior to medication initiation are an example. We chose to model these hospitalizations in three ways: multiple indicators for the overall number of recent nonmental health hospitalizations, whether any hospitalization of a particular type (e.g. ICU, cardiac, etc.) had occurred in the past 2 years, and what type of hospitalization had occurred most recently prior to medication initiation. However, it can be easily conceptualized that near-complete modeling of hospitalizations experienced by the patient in the last two years might have included the timing and number of days preceding medication initiation for every hospitalization type, and potentially length of stay as well. Furthermore, indicators concerning whether multiple hospitalizations of the same type had occurred and how separated in time the repeat admissions were might be desirable. And of course, as pointed out in the manuscript, we did not have information about hospitalizations that occurred outside the VHA system. For many variables, such as those denoting recent hospitalizations, at some point practical decisions must be made concerning what detail in modeling is appropriate and feasible.

### Online Supplement DS13  The Importance of Identifying Incident Users

Including prevalent and/or past users of a medication in a study cohort can induce a number of potential undesirable biases.[30] For instance, including prevalent users can introduce "survivor bias" and/or "adherence bias" (sometimes identified by the overarching term "healthy users" or "healthy adherers",[84] complicate the control of confounding (since "baseline" covariate values may be affected by the treatment itself), and select for patients who find the medication either unusually tolerable or effective. Including past users can also introduce bias since patients who have had a past trial of a medication are likely to select it again if they (or their provider) believe it to have been effective, and are likely to choose an alternative if they (or their provider) believe it to have likely been ineffective.

However, a design choice is forced in nonrandomized research concerning how long of a "clean period" to implement as an exclusion criteria for cohort membership. What is usually done is to impose some uniform requirement, such as 12 months with no evidence of a prescription of either study medication. Obviously, such a requirement creates the possibility that "past users" with use more remote than 12 months will be included in the study cohort. Despite such a drawback, such a uniform period is usually chosen for nonrandomized studies comparing medications, because the alternative approach, prohibiting any past use, risks creating a variably-determined "clean period" between patients. If any past use is prohibited, then longtime users of a system are forced to not be exposed to the study medication for 3 years, 5 years, 7 years, etc. (depending on the length of their electronic medical record), whereas patients more recently entering the cohort by default do not have such a restriction, since their medical record does not extend back as far.

In this study, because we had some concern that we might have inadequate sample size to examine a uniform follow-up period (i.e., 365 days from initiation), we deliberately chose a less restrictive "clean period" of 6 months. However, for the small number of patients for which we had evidence of past use of either medication, a covariate was included in the high-dimensional propensity score denoting past use of either medication. (This approach is the approach recommended when the "incident user" design was described[30]). This approach led to an extremely similar proportion of patients in both treatment groups having prior exposure to one of the two medications (11.9% of the hdPS-patients initiating valproate and 12.0% of the hdPS-matched patients initiating lithium). While these procedures do not eliminate entirely the possibility for bias, especially in VHA studies in which outside medication use is a possibility, they likely reduce substantially the impacts of any bias relating to prevalent or past users.

**Additional references**

77 Toh S, Garcia Rodriguez LA, Hernan MA. Confounding adjustment via a semi-automated high-dimensional propensity score algorithm: an application to electronic medical records. Pharmacoepidemiol Drug Saf. 2011; **20**(8): 849-57.
78 Shack LG, Rachet B, Williams EM, Northover JM, Coleman MP. Does the timing of comorbidity affect colorectal cancer survival? A population based study. Postgrad Med J. 2010; **86**(1012): 73-8.
79 Quan H, Sundararajan V, Halfon P, Fong A, Burnand B, Luthi JC, et al. Coding algorithms for defining comorbidities in ICD-9-CM and ICD-10 administrative data. Med Care. 2005; **43**(11): 1130-9.
80 Haukka J, Suominen K, Partonen T, Lonnqvist J. Determinants and outcomes of serious attempted suicide: a nationwide study in Finland, 1996-2003. Am J Epidemiol. 2008; **167**(10): 1155-63.
81 Kapur N, Cooper J, King-Hele S, Webb R, Lawlor M, Rodway C, et al. The repetition of suicidal behavior: a multicenter cohort study. J Clin Psychiatry. 2006; **67**(10): 1599-609.
82 Kim HM, Smith EG, Stano CM, Ganoczy D, Zivin K, Walters H, et al. Validation of key behaviourally based mental health diagnoses in administrative data: suicide attempt, alcohol abuse, illicit drug abuse and tobacco use. BMC Health Serv Res. 2012; **12**: 18.
83 Centers for Disease Control. National Center for Injury Prevention and Control. WISQARS Injury Mortality Reports, 1999-2007. [cited June 2, 2012. ]; Available from:
http://webappa.cdc.gov/sasweb/ncipc/mortrate10_sy.html
84 Simpson SH, Eurich DT, Majumdar SR, Padwal RS, Tsuyuki RT, Varney J, et al. A meta-analysis of the association between adherence to drug therapy and mortality. BMJ. 2006; **333**(7557): 15.