Take a deep dive into the business of transplantation!

# Liver Transplant Financial Bootcamp
# &
# Kidney Transplant Financial Bootcamp

*Co-located for 2019!*

September 19–20, 2019

Loews Chicago O'Hare Hotel

**Register at ASTS.org/bootcamps**

**ASTS**
American Society of Transplant Surgeons®

# Big Data, Predictive Analytics, and Quality Improvement in Kidney Transplantation: A Proof of Concept

T. R. Srinivas[1,*,†], D. J. Taber[2], Z. Su[3], J. Zhang[3], G. Mour[1], D. Northrup[4], A. Tripathi[5], J. E. Marsden[3], W. P. Moran[3] and P. D. Mauldin[3]

[1]Division of Nephrology, Medical University of South Carolina, Charleston, SC
[2]Division of Transplant Surgery, Medical University of South Carolina, Charleston, SC
[3]Division of General Internal Medicine & Geriatrics, Medical University of South Carolina, Charleston, SC
[4]Office of the Chief Information Officer, Medical University of South Carolina, Charleston, SC
[5]IBM Corporation, Armonk, NY
*Corresponding author: Titte R. Srinivas,
titte.srinivas@imail.org
[†]Current address: Intermountain Transplant Services, Salt Lake City, UT

We sought proof of concept of a Big Data Solution incorporating longitudinal structured and unstructured patient-level data from electronic health records (EHR) to predict graft loss (GL) and mortality. For a quality improvement initiative, GL and mortality prediction models were constructed using baseline and follow-up data (0–90 days posttransplant; structured and unstructured for 1-year models; data up to 1 year for 3-year models) on adult solitary kidney transplant recipients transplanted during 2007–2015 as follows: Model 1: United Network for Organ Sharing (UNOS) data; Model 2: UNOS & Transplant Database (Tx Database) data; Model 3: UNOS, Tx Database & EHR comorbidity data; and Model 4: UNOS, Tx Database, EHR data, Posttransplant trajectory data, and unstructured data. A 10% 3-year GL rate was observed among 891 patients (2007–2015). Layering of data sources improved model performance; Model 1: area under the curve (AUC), 0.66; (95% confidence interval [CI]: 0.60, 0.72); Model 2: AUC, 0.68; (95% CI: 0.61–0.74); Model 3: AUC, 0.72; (95% CI: 0.66–077); Model 4: AUC, 0.84, (95 % CI: 0.79–0.89). One-year GL (AUC, 0.87; Model 4) and 3-year mortality (AUC, 0.84; Model 4) models performed similarly. A Big Data approach significantly adds efficacy to GL and mortality prediction models and is EHR deployable to optimize outcomes.

Abbreviations: AUC-ROC, area under the curve–receiver operating characteristic curve; BK, BK virus; CI, confidence interval; CMV, cytomegalovirus; DGF, delayed graft function; eGFR, estimated glomerular filtration rate; EHR, electronic health record; GL, graft loss; ICD-9, international classification of diseases; KDRI, kidney donor risk index; MI, myocardial infarction; NLP, natural language processing; OR, odds ratio; PCR, polymerase chain reaction; SRTR, Scientific Registry of Transplant Recipients; Tx database, transplant database; UNOS, United Network for Organ Sharing

## Introduction

Despite continuing improvement in short-term graft outcomes, long-term graft survival has not substantially improved (1). This disconnect exists in a wealth of readily available outcomes data in large national data registries such as those maintained by the United Network for Organ Sharing (UNOS) and the Scientific Registry of Transplant Recipients (SRTR) (2). While the SRTR and UNOS capture a large amount of baseline information and data relating to graft loss (GL) and death with remarkable completeness, they are not built to capture longitudinally evolving patient-level data (2). This deficiency limits predictive accuracy with regard to GL and mortality, and their utility in quality improvement initiatives or data-driven clinical decision making. With the explosion in the amounts of patient-level data available through electronic health records (EHRs), analyses could potentially incorporate dynamically evolving clinically relevant patient-level data in a meaningful manner that may inform on the care of individual patients or subpopulations of patients at highest risk of GL or death. Using manually extracted EHR data, we have previously demonstrated that addition of dynamic patient-level data improves predictive accuracy over administrative data for 30-day readmissions following kidney transplantation (demonstrated by others as a surrogate for subsequent GL) (3,4). However, broad applicability of such an approach outside research contexts is limited by practical constraints to achieving near real-time capture of structured clinical data and more importantly, the need for manual abstraction of unstructured data fields (e.g. clinical notes, pathology reports) to populate analyses.

Patient-level data available through EHRs are notable for their volume, velocity, and variety, all of which define the scope of Big Data (5). Data in the EHR are also notable for information residing in free text, deemed "unstructured data" in that they cannot be readily applied in analyses without being suitably abstracted. Natural language processing (NLP) describes computational techniques that explore the interaction between computers and human (natural) language. This technique can be used to extract information from text fields and allow their incorporation into analyses without the need for manual abstraction.

The recent availability of high-throughput Big Data approaches to collect and curate structured and unstructured data coupled with high-throughput integrated analytic solutions affords the promise of incorporating large data sets in analyses in near real-time fashion and thereby bringing predictive models with improved accuracy to the bedside to drive preventive or therapeutic intervention (5). In the context of kidney transplantation, structured laboratory data such as serum creatinines and glomerular filtration rates (GFRs) and hemoglobin values are obtained in large numbers over time, as are unstructured data elements buried in text form (e.g. biopsy reports, dictated vital signs, social worker notes). As such, volume, velocity, and variety of data that constitute a definition of Big Data would be typified by kidney transplant patient-level data. We began with the premise that this abundance of patient-level data affords the ideal stage to prove the concept that incorporating longitudinally evolving structured and unstructured patient-level data using Big Data approaches could improve our ability to predict GL and mortality among kidney transplant recipients compared to what is possible with national data.

We report on a multidisciplinary quality improvement initiative that began with the hypothesis that prediction of 1- and 3-year GL and mortality among kidney transplant recipients could be substantially improved using dynamic models that utilize baseline clinical data and patient-level information acquired as a component of clinical care, as compared to predictive models derived using only structured static variables.

## Methods

### Population and setting
This is a single-center retrospective cohort quality improvement initiative among adult solitary kidney transplant recipients ≥18 years of age transplanted at the Medical University of South Carolina, Charleston, SC, between January 2007 and June 2015. We chose this period as protocols for patient and regimen selection were better standardized and data elements within the EHR were better defined and collected. Patients were excluded if they (1) had a GL or death in the first 7 days posttransplant, or (2) did not have a GL or death recorded nor any other data during the first year posttransplant (for the 1-year model) or first 3 years posttransplant (for the 3-year model). The primary immunosuppressive protocol comprised tacrolimus, mycophenolate mofetil, and corticosteroids with IL-2 receptor antagonist or rabbit antithymocyte globulin induction.

### Data sources
Detailed structured data were directly acquired from electronic medical records (EHR); (Practice Partner [Pre-May 2012] [McKesson, Seattle, WA]); Epic July 2011 onwards (Epic Corporation, Madison, WI), UNOS database elements containing Organ Procurement and Transplantation Network data since 1986. Our Transplant Database (Tx Database) (Velos, Inc., Fremont, CA) was used to extract key social determinants of health (Supplemental File). NLP was applied to unstructured text fields using proprietary NLP solutions in the IBM Watson Content Analytics (IBM Corporation, Armonk, NY) to extract Banff scores and vital signs that predated electronic capture (6). Using NLP algorithms, we parsed Banff lesion scores from pathology reports in text form. Lesion scores transcribed as g0, t0, i2, t2, v0 were extracted and transferred to analytic databases (as an example, g = 0, t = 0, i = 2, and t = 2 would constitute Pathologic grade, Type IIa rejection). Furthermore, values transcribed in error as "tII" were deemed semantically equivalent to t2 and so extracted and analyzed.

### Primary outcome measures
GL and mortality information were retrieved from Tx Database and UNOS files. Account of any GL within 1 year or 3 years was defined as a return to chronic dialysis, retransplantation, or death. The unit of observation being the unique transplant, multiple GL per patient were considered unique observations. If patients received more than one transplant, death was linked to the most recent transplant event. A 90-day exposure period was used for 1-year GL and mortality models and a 1-year exposure period was used to derive 3-year GL and mortality models.

### Covariates
UNOS data elements were utilized for key demographic and transplant-related variables in accordance with published methodology used by the SRTR (see supplemental file). The Tx database was to extract key social determinants of health (see supplemental file). EHR data were utilized to supplement obesity data and vital signs, as well as to provide comorbidities, cardiovascular events, laboratory data, transplant length of stay, and posttransplant acute care utilization data, both inpatient and emergency department.

Means, standard deviations, maximums, and regressed slopes were used to represent dynamic variables, capturing effects of change, direction of change, and magnitude of change throughout the exposure period. This treatment was applied to estimated GFRs (eGFRs), pulse rates, blood pressures, and hemoglobin levels (see supplemental file).

Using enhanced International Classification of Diseases (ICD-9-CM) codes, comorbidity was derived from a modified Elixhauser coding algorithm and select Charlson comorbidities (7). Event data were captured beginning at transplant date and extending throughout the exposure period. In the conceptual model build, we used severable variables as surrogates to capture the many risk domains for GL and mortality (8). For instance, kidney donor risk index (KDRI) and evolution of blood pressures were used as surrogates for kidney quality. Hemoglobin slopes, eGFR evolution, and delayed graft function (DGF) were physiologic surrogates for graft function. Posttransplant cardiovascular events such as arrhythmias and myocardial infarction reflected cardiovascular risk. Immunologic risk was captured by rejection rates (Banff scores), cytomegalovirus (CMV) infection, BK virus (BK) infection, and tacrolimus trough concentrations. Social determinants of health were captured in demographics, caregiver status, education level, and income. Unplanned acute care utilization reflected evolving acuity and access to care (Supplemental Table).

## Statistical analysis

The goal of this quality improvement initiative was to develop a prediction model that derives a patient-level risk score that can be used at the bedside that incorporates "real-time" data interactions to determine potential changes in a patient's probability of GL by 1 and 3 years, and mortality by 3 years (where the timing of the event within those time periods is not a priority). Three predictive risk models were developed, using baseline and follow-up data (up to 90 days posttransplant exposure period for the 1-year model [1 year GL]), up to 365 days posttransplant exposure period for the 3-year models (3 years GL and 3 years mortality), from both structured and unstructured data formats. The Firth (to account for a low number of events (9–11)) multivariable logistic regression was employed. Statistical significance was determined at the two-sided 5% level. A combination of statistical and clinical information was used for variable selection. First, the backwards selection process in logistic Firth using an exit p-value at the 20% level was used. Next, variables were chosen in parallel based on information from the stepwise AIC variable selection criteria (Data S1). Clinical adjudication was used if discrepancies of variables were revealed between methods. To assess and adjust for potential model overfitting, we used the Harrell Optimism Correction (Data S1) (12). Once the final model was selected, Bootstrapping methodology (1000 iterations) was used for model internal validation and area under the curve (AUC) was used to determine and compare model accuracy. To assess the potential for biases due to noncensored data, we ran survival models using proportional hazard methodology (results reported in supplementary file). Finally, an empirical cross check between actual observed events and mod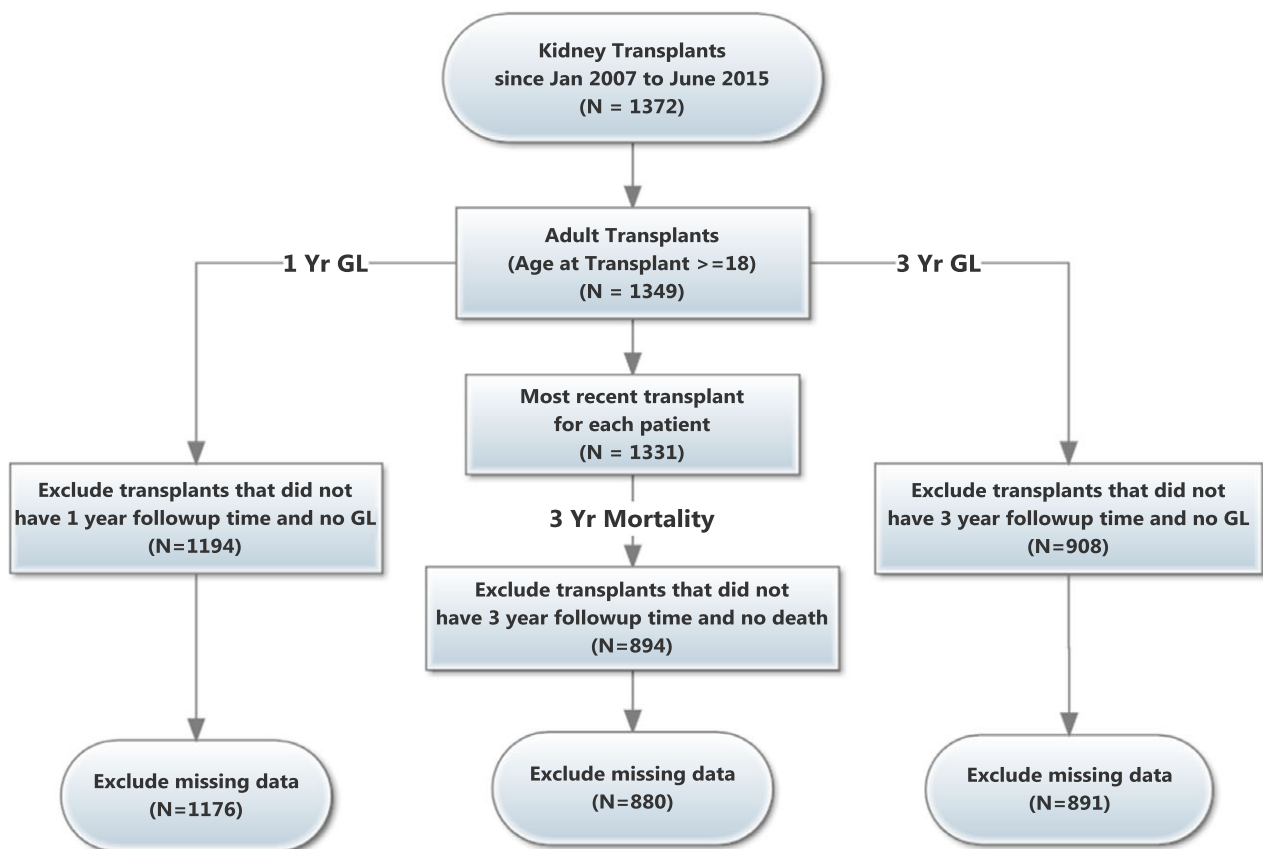el predicted risk scores was performed by clinical experts. IBM Watson Content Analytics Suite IBM SPSS Modeler (Version 17) and IBM SPSS Statistics - Essentials for R (IBM Corporation) were used for the statistical analysis.

## Results

Of 1372 kidney transplants since January 2007, 1176 transplant patients met eligibility criteria for the 1-year GL model, 891 for the 3-year GL model, and 880 for the 3-year mortality model (Figure 1; Table 1: Demographic and clinical characteristics of the study population, stratified by event). Statistically significant differences existed for a number of variables, as summarized below and in Table 2.

### Donor quality, blood pressures, and pulse rates

A 1-unit increase in donor's KDRI score increased the patient's likelihood of a 3-year GL by over three times (3.06 odds ratio [OR]; [1.49, 6.28] 95% CI). Greater variability in pulse pressures was associated with significantly higher GL odds at 1 and 3 years (1.14 OR each; [1.04, 1.25], and [1.06, 1.22] 95% CI, respectively), while controlling for systolic blood pressure. A 1-unit increase in mean pulse rate over the first year posttransplant was associated with a 3% increase in odds of 3-year GL (1.03 OR;



**Figure 1: Description of patients included or excluded in this report.** GL, graft loss.

**Table 1:** Comparison of patient characteristics

| | 1-year graft loss (GL) N(GL) = 45\|N(No GL) = 1131 | 3-year graft loss (GL) N(GL) = 89\|N(No GL) = 802 | 3-year mortality N(Dead) = 76\| N(Alive) = 804 |
|---|---|---|---|
| UNOS | | | |
| Age at transplant | 54 ± 14\|52 ± 14 | 52 ± 14\|52 ± 14 | 58 ± 13\|52 ± 14* |
| Female | 31%\|41% | 30%\|42%* | 36%\|42% |
| African American | 58%\|53% | 56%\|53% | 47%\|54% |
| African American donor | 29%\|26% | 30%\|27% | 28%\|27% |
| KDRI | 1.2 ± 0.4\|1.0 ± 0.3* | 1.1 ± 0.4\|1.0 ± 0.3* | 1.1 ± 0.4\|1.0 ± 0.4* |
| Blood type B | 29%\|16%* | 23%\|16% | 16%\|16% |
| Waitlisting time (years) | 1.5 ± 1.6\|1.7 ± 1.5 | 1.8 ± 1.6\|1.6 ± 1.4 | 1.9 ± 1.8\|1.6 ± 1.4 |
| First week dialysis | 27%\|13%* | 15%\|8%* | 24%\|7%* |
| Diabetes | 27%\|34% | 30%\|32% | 45%\|32%* |
| Obesity (BMI ≥30) | 31%\|42% | 37%\|40% | 36%\|40% |
| Private insurance | 24%\|28% | 21%\|29% | 15%\|29%* |
| Distance to MUSC (miles) | 159 ± 95\|155 ± 123 | 156 ± 98\|156 ± 122 | 155 ± 103\|158 ± 125 |
| Immobility (functional status ≤40%) | 0%\|0.6% | 0%\|0.7% | 1.3%\|0.7% |
| Previous kidney transplant | 11%\|8.8% | 9%\|8.6% | 15%\|8% |
| Graft loss by 1 year | | | 17%\|3%* |
| Velos | | | |
| Finish high school | 62%\|76%* | 66%\|74% | 71%\|73% |
| Employed | 16%\|29% | 23%\|29% | 21%\|28% |
| Received disability | 47%\|42% | 43%\|40% | 37%\|41% |
| Married | 47%\|54% | 45%\|65% | 55%\|52% |
| Smoker | 16%\|7%* | 10%\|8% | 9%\|8% |
| Primary caregiver | 76%\|79% | 63%\|78%* | 72%\|77% |
| EHR | | | |
| Congestive heart failure | 18%\|12% | 16%\|12% | 20%\|12%* |
| Peripheral vascular disorders | 11%\|10% | 9%\|9% | 16%\|9%* |
| Cerebrovascular disease | 0%\|5% | 1%\|6% | 3%\|6% |
| Cardiac arrhythmias | 42%\|21%* | 32%\|21%* | 36%\|21%* |
| Valvular disease | 16%\|9% | 11%\|10% | 16%\|9% |
| Hypertension | 96%\|96% | 97%\|96% | 99%\|96% |
| Alcohol abuse | 9%\|3%* | 6%\|3% | 4%\|3% |
| Drug abuse | 7%\|2% | 5%\|2% | 1%\|2% |
| Depression | 9%\|12% | 16%\|10% | 7%\|11% |
| Myocardial infarction | 13%\|7% | 10%\|7% | 15%\|7%* |
| Transplant LOS (days) | 4.3 ± 3.2\|3.5 ± 2.2 | 3.8 ± 2.4\|3.4 ± 3.0 | 4.9 ± 6.6\|3.4 ± 1.2* |
| Acute myocardial infarction during exposure[1] | 2.2%\|0.6% | 4.5%\|0.6%* | 3.9%\|0.6%* |
| Cardiac or vascular event during exposure | 40%\|13%* | 44%\|18%* | 50%\|18%* |
| BK>500 during exposure | 4%\|5% | 11%\|10% | 16%\|10% |
| CMV>500 during exposure | 2%\|8% | 17%\|18% | 22%\|17% |
| Pulse mean during exposure | 81 ± 10\|78 ± 9 | 80 ± 9\|77 ± 8* | 80 ± 10\|77 ± 8* |
| SBP mean during exposure | 137 ± 16\|142 ± 14* | 140 ± 14\|140 ± 12 | 138 ± 14\|140 ± 13 |
| Pulse pressure SD during exposure | 14.7 ± 4.1\|13.5 ± 3.9* | 15.4 ± 4.1\|13.9 ± 3.8* | 15.0 ± 4.0\|14.0 ± 3.9* |
| Glucose mean during exposure | 128 ± 28\|135 ± 32 | 138 ± 33\|133 ± 32 | 141 ± 28\|133 ± 32 |
| HGB mean – day 7 until end of exposure | 9.9 ± 1.5\|10.9 ± 1.4* | 10.5 ± 1.6\|11.3 ± 1.4* | 10.4 ± 1.5\|11.3 ± 1.4* |
| HGB slope per month – day 7 until end of exposure | −1.5 ± 8.1\|0.56 ± 0.93 | −0.72 ± 5.76\|0.32 ± 0.63 | 0.15 ± 0.73\|0.31 ± 0.61* |
| Max eGFR during exposure | 47 ± 30\|66 ± 24* | 61 ± 47\|74 ± 26* | 75 ± 55\|73 ± 28 |
| eGFR SD during exposure | 10.3 ± 7.0\|16.1 ± 6.9* | 13.6 ± 9.1\|16.5 ± 6.8* | 16.2 ± 10.8\|16.4 ± 6.9 |
| eGFR slope from max value until end of exposure | −0.8 ± 1.5\|−0.5 ± 1.4 | 0.5 ± 1.0\|−0.3 ± 1.0 | −0.8 ± 2.8\|−0.2 ± 0.7 |
| Tacrolimus mean during exposure | 8.6 ± 1.4\|9.0 ± 1.4 | 8.4 ± 1.5\|8.6 ± 1.4 | 8.6 ± 1.5\|8.7 ± 1.3 |
| Tacrolimus SD during exposure | 3.9 ± 1.2\|3.9 ± 1.3 | 4.1 ± 1.4\|3.9 ± 1.3 | 4.2 ± 1.3\|3.9 ± 1.3* |
| Days since TX to first tacrolimus 8–12 during exposure | 6 ± 5\|7 ± 7 | 6 ± 4\|8 ± 9 | 7 ± 11\|7 ± 9 |
| Inpatient readm count during exposure | 1.2 ± 1.5\|0.5 ± 0.8* | 1.7 ± 1.9\|0.8 ± 1.2* | 2.0 ± 2.0\|0.8 ± 1.3* |

(*Continued*)

**Table 1.** Continued

| | 1-year graft loss (GL) N(GL) = 45\|N(No GL) = 1131 | 3-year graft loss (GL) N(GL) = 89\|N(No GL) = 802 | 3-year mortality N(Dead) = 76\| N(Alive) = 804 |
|---|---|---|---|
| Emergency department visit count during exposure | 0.04 ± 0.21\|0.05 ± 0.26 | 0.2 ± 1.0\|0.1 ± 0.5 | 0.3 ± 0.6\|0.1 ± 0.6 |
| NLP | | | |
| Maximum acute Banff score during exposure | 1.4 ± 2.1\|0.5 ± 1.3* | 2.3 ± 2.6\|0.7 ± 1.5* | 1.2 ± 1.6\|0.9 ± 1.7 |

BK, BK virus; CMV, cytomegalovirus; EHR, electronic health record; eGFR, estimated glomerular filtration rate; HGB, hemoglobin; KDRI, kidney disease risk index; LOS, length of stay; MUSC, Medical University of South Carolina; NLP, natural language processing; SBP, systolic blood pressure; SD, standard deviation; TX, transplant; UNOS, United Network for Organ Sharing.
[1]Exposure period for 1-year model is 90 days. Exposure period for 3-year models is 365 days.
*p-value ≤0.05.

[1.00, 1.07] 95% CI) and a 4% increase in the odds of 3-year mortality (1.04 OR; [1.00, 1.07] 95% CI) (Table 2).

### Posttransplant clinical evolution of graft function
Patients with DGF experienced significantly higher odds for 3-year mortality (2.48 OR; [1.20, 5.00] 95% CI), as compared to those without DGF. Those patients whose GFR continued to improve in the first 90 days were at significantly lower risk for graft loss. Higher maximum values of eGFR labs over the exposure periods were associated with significant reduction in 1- and 3-year GL rates (0.97 and 0.99 OR; [0.95, 0.99] and [0.98, 1.00] 95% CI, respectively). An increased eGFR slope from the maximum value over the first year posttransplant to day 365 posttransplant reduced the odds of 3-year mortality by 13% (0.87 OR; [0.73, 0.98] 95% CI), and an increased eGFR slope from the maximum value over the first 90 days posttransplant to day 90 posttransplant also reduced the odds of 1-year GL by 17% (0.83 OR; [0.74, 0.99] 95% CI). An increasing hemoglobin slope from day 7 posttransplant to end of exposure period was associated with a substantial decrease in odds for both 1- and 3-year GL (0.78 and 0.72 OR; [0.58, 0.93] and [0.56, 0.87] 95% CI, respectively). Finally, the presence of a 1-year GL increased the odds of 3-year mortality by nearly threefold (2.77 OR; [1.13, 6.54] 95% CI).

### Demographic and waitlist characteristics
Increasing recipient age was associated with reduced relative odds for 3-year GL, and a slight increased odds for 3-year mortality (0.98 and 1.03 OR; [0.96, 1.00] and [1.00, 1.05] 95% CI, respectively). African American race was associated with a reduced relative odds of 3-year mortality compared to other races (0.46 OR; [0.26, 0.81] 95% CI). Females were >40% less likely than males to have a 3-year GL (0.57 OR; [0.32, 0.99] 95% CI). Increased transplant waitlisting times (years) were associated with lower odds of 1-year GL and patients with blood type B had over three times the likelihood of having a 1-year GL compared to those with other blood types (3.41 OR; [1.55, 7.35] 95% CI).

### Associations of immunologic risk
Acute rejection at 1 year (acute Banff lesion scores) was associated with a higher 3-year GL (1.37 OR; [1.22, 1.54] 95% CI). Those with at least one positive cytomegalovirus polymerase chain reaction (PCR) (CMV-PCR) copy number over 500 by 90 days posttransplant showed 75% lower relative odds of 1-year GL (0.20 OR; [0.02, 0.88] 95% CI) than those with no positive CMV-PCR. There was a trend for those with a positive BK PCR in the first year to have almost double the odds of 3-year mortality (1.93 OR; [0.90, 3.88] 95% CI).

### Social determinants of health
Patients with private insurance showed lower relative odds of 3-year mortality than those with nonprivate insurance (0.46 OR; [0.21, 0.91] 95% CI), and those patients who completed high school were more than 50% less likely to have a 1-year GL (0.47 OR; [0.23, 0.97] 95% CI). Patients who identified a primary caregiver at the time of transplant were significantly less likely to have a 3-year GL (0.40 OR; [0.23, 0.69] 95% CI).

### Cardiovascular risk factors
Pretransplant cardiovascular comorbidity was associated with posttransplant cardiovascular events. Patients experiencing acute myocardial infarctions over the first year following transplant were over 11 times more likely to experience a 3-year GL (11.14 OR; [2.15, 54.30] 95% CI) than those without the event. Similarly, patients with cardiac or vascular events over the associated exposure periods were at significant increased risk for poor outcomes: nearly 2.5 times more likely to have a 1-year GL (2.48 OR; [1.06, 5.66] 95% CI), nearly three times more likely for a 3-year GL (2.98 OR; [1.74, 5.10] 95% CI), and nearly 2.25 times more likely to experience 3-year mortality (2.23 OR; [1.21, 4.08] 95% CI).

### Acute care utilization
Each additional day of length of stay associated with the transplant hospitalization associated with higher relative odds of 3-year mortality (1.08 OR; [1.01, 1.26] 95% CI)

**Table 2:** Multivariable logistic (Firth) regression odds ratio and 95% confidence intervals (CI) (Model 4: UNOS + Velos + EHR comorbidity + EHR posttransplant trajectory + NLP)

| | Odds ratio (95% CI) | | |
| --- | --- | --- | --- |
| | 1-year graft loss | 3-year graft loss | 3-year mortality |
| UNOS | | | |
|   Age at transplant | | 0.98 (0.96, 1.00)* | 1.03 (1.00, 1.05)* |
|   Female | | 0.57 (0.32, 0.99)* | |
|   African American | | | 0.46 (0.26, 0.81)* |
|   KDRI | 2.48 (0.95, 6.43) | 3.06 (1.49, 6.28)* | 1.81 (0.88, 3.69) |
|   Blood type B | 3.41 (1.55, 7.35)* | | |
|   Waitlisting time (years) | 0.76 (0.56, 0.98)* | | |
|   First week dialysis | | | 2.48 (1.20, 5.00)* |
|   Diabetes | | | 1.67 (0.93, 3.02) |
|   Obesity | 0.50 (0.23, 1.06) | 0.66 (0.38, 1.11) | 0.65 (0.36, 1.14) |
|   Private insurance | | | 0.46 (0.21, 0.91)* |
|   Previous kidney transplant | | | 2.31 (0.99, 5.05) |
|   Graft loss by 1 year | | | 2.77 (1.13, 6.54)* |
| Transplant database | | | |
|   Finish high school | 0.47 (0.23, 0.97)* | | |
|   Smoker | 2.61 (0.87, 6.91) | | |
|   Primary caregiver identified at transplant | | 0.40 (0.23, 0.69)* | |
| EHR | | | |
|   Cerebrovascular disease | 0.07 (<0.01, 0.65)* | 0.23 (0.02, 1.13) | |
|   Cardiac arrhythmias | 2.16 (1.01, 4.52)* | | |
|   Alcohol abuse | 3.71 (0.87, 12.73) | 3.22 (0.89, 9.78) | |
|   Drug abuse | 3.55 (0.63, 15.50) | | |
|   Depression | | 1.91 (0.86, 3.98) | 0.44 (0.14, 1.13) |
|   Transplant LOS (days) | | | 1.08 (1.01, 1.26)* |
|   Acute MI during exposure[1] | | 11.14 (2.15, 54.30)* | |
|   Cardiac or vascular event during exposure | 2.48 (1.06, 5.66)* | 2.98 (1.74, 5.10)* | 2.23 (1.21, 4.08)* |
|   BK>500 during exposure | | | 1.93 (0.90, 3.88) |
|   CMV>500 during exposure | 0.20 (0.02, 0.88)* | | |
|   Pulse mean during exposure | | 1.03 (1.00, 1.07)* | 1.04 (1.00, 1.07)* |
|   SBP mean during exposure | 0.97 (0.94, 1.00)* | | |
|   Pulse pressure SD during exposure | 1.14 (1.04, 1.25)* | 1.14 (1.06, 1.22)* | |
|   Glucose mean during exposure | 0.99 (0.98, 1.00) | | |
|   HGB slope per month: day 7 until end of exposure | 0.78 (0.58, 0.93)* | 0.72 (0.56, 0.87)* | |
|   Max eGFR during exposure | 0.97 (0.95, 0.99)* | 0.99 (0.98, 1.00)* | |
|   eGFR slope from max value until end of exposure | 0.83 (0.74, 0.99)* | | 0.87 (0.73, 0.98)* |
|   Tacrolimus SD during exposure | | | 1.18 (0.97, 1.42) |
|   Inpatient readm count during exposure | 1.42 (1.03, 1.93)* | | 1.16 (0.98, 1.37) |
|   Emergency department visit count during exposure | | 1.31 (0.93, 1.74) | |
| NLP | | | |
|   Max acute Banff score during exposure | | 1.37 (1.22, 1.54)* | |

BK, BK virus; CMV, cytomegalovirus; EHR, electronic health record; eGFR, estimated glomerular filtration rate; HGB, hemoglobin; KDRI, kidney disease risk index; LOS, length of stay; MI, myocardial infarction; NLP, natural language processing; readm, readmission; SBP, systolic blood pressure; UNOS, United Network for Organ Sharing.
[1]Exposure period for 1-year model is 90 days. Exposure period for 3-year models is 365 days.
*p-value ≤0.05.

as did higher rates of inpatient readmission in the first 90 days posttransplant (1.42 OR; [1.03, 1.93] 95% CI).

### Layering of data sources augments predictive accuracy

Table 3 and Figure 2 demonstrate the added value of data sources and variable construction on the accuracy of the iterative predictive models. For the 1-year GL model, if only UNOS data were used in the predictive model, the AUC-ROC was 0.716 (0.641, 0.790 95% CI; Model 1). With the addition of caregiver data from the transplant database, the predictive performance improved, with an AUC of 0.741 (0.669, 0.814 95% CI; Model 2). EHR-derived comorbidity data improved the accuracy of the models' predictive capability, with an AUC of 0.769 (0.692, 0.845 95% CI; Model 3). With the addition of trajectory and NLP variables to the model, the AUC significantly improved to 0.873 (0.807, 0.939, 95%

**Table 3:** Comparison of AUC-ROC curves between three primary models using varying data sources

| 1-year graft loss | | 3-year graft loss | | 3-year mortality | |
|---|---|---|---|---|---|
| Data model A vs Data model B | p-value | Data model A vs Data model B | p-value | Data model A vs Data model B | p-value |
| 2 (AUC = 0.741) vs 1 (AUC = 0.716) | 0.416 | 2 (AUC = 0.665) vs 1 (AUC = 0.661) | 0.838 | 2 (AUC = 0.765) vs 1 (AUC = 0.765) | 1.000 |
| 3 (AUC = 0.769) vs 2 (AUC = 0.741) | 0.119 | 3 (AUC = 0.712) vs 2 (AUC = 0.665) | 0.006 | 3 (AUC = 0.768) vs 2 (AUC = 0.765) | 0.476 |
| 4 (AUC = 0.873) vs 3 (AUC = 0.769) | 0.001 | 4 (AUC = 0.846) vs 3 (AUC = 0.712) | <0.001 | 4 (AUC = 0.838) vs 3 (AUC = 0.768) | 0.007 |

AUC-ROC, UC, area under the curve–receiver operating characteristic curve; EHR, electronic health record; NLP, natural language processing; UNOS, United Network for Organ Sharing.
Data Model 1: UNOS only.
Data Model 2: UNOS + Transplant database.
Data Model 3: UNOS + Transplant database + EHR comorbidity.
Data Model 4: UNOS + Transplant database + EHR comorbidity + EHR posttransplant trajectory + NLP.

CI; Model 4). Similar improvements in iterative model accuracy with layering of data sources were demonstrated for the outcomes of 3-year GL and 3-year mortality, with significant increases in AUCs for 3-year GL model ranging from 0.661 ([0.598, 0.724] 95% CI: Model 1) to 0.846 ([0.798, 0.894] 95% CI: Model 4), and 3-year mortality increasing from AUCs of 0.765 ([0.702, 0.827] 95% CI: Model 1) to 0.838 ([0.790, 0.885] 95% CI: Model 4).
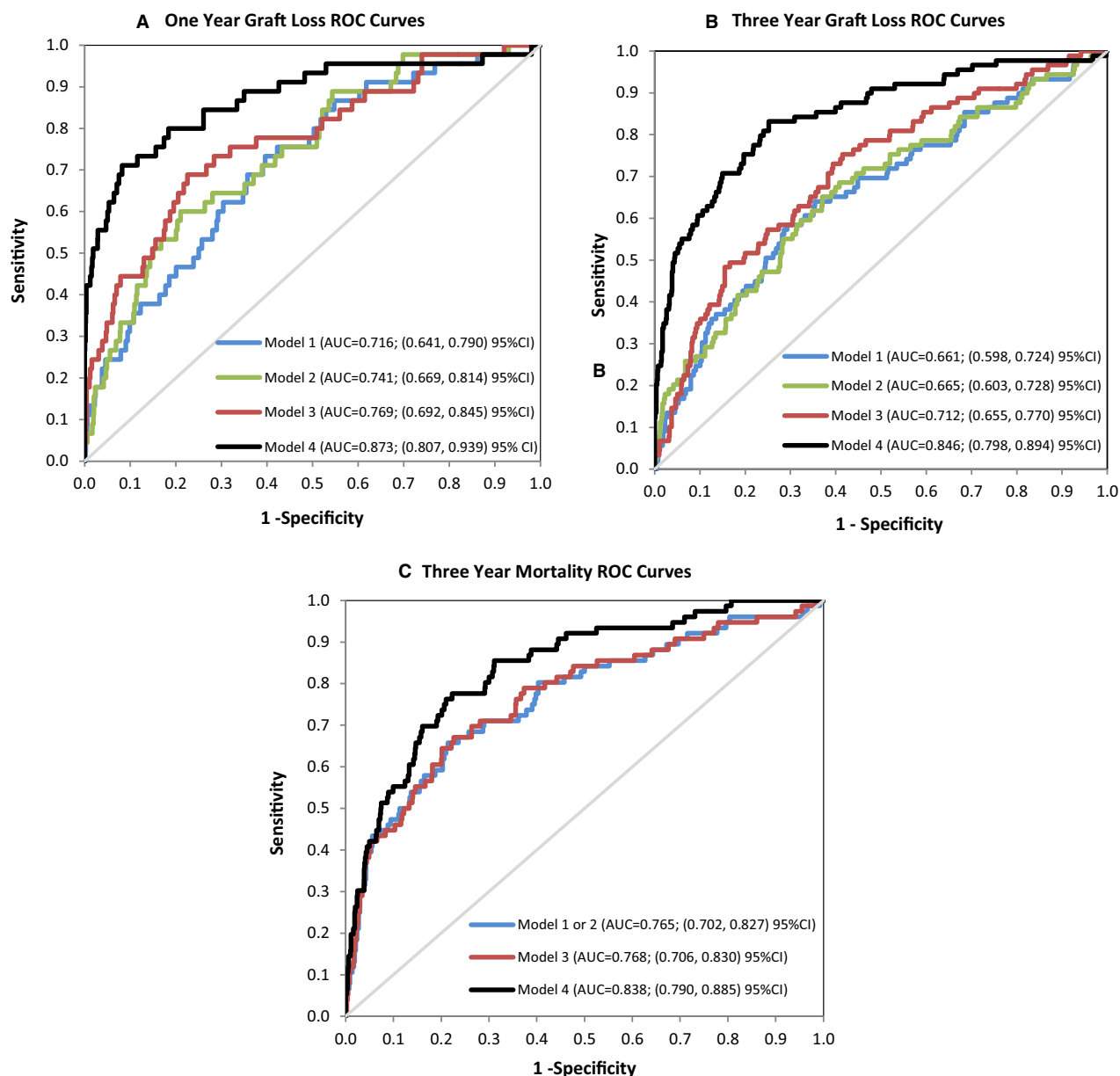
### *Model checking*

Internal validation was conducted for the final models (Model 4 for each of the three outcomes). Bootstrapping demonstrated these predictive models to be stable, with minimal bias or deviations within the standard errors and optimism-corrected AUCs remained stable (values in Data S1). Finally, the Cox proportional hazard models demonstrated similar results with the logistic models (results provided in Data S1).

## Discussion

We have demonstrated proof of concept of an approach to modeling transplant outcomes using data that are available in the EHR that substantially improve the accuracy of predictive models for GL and mortality. The current predictive ability of the national models is modest, with the most recently published c-statistic being 0.68 (13).

Our results reflect the dominant mechanistic underpinnings of the biology of transplantation and how they relate to the environment of care that the transplanted organ and the patient interact with, in turn, discernible as clinical outcomes. The significance of higher KDRI as a strong correlate of GL reflects kidney quality; higher kidney quality associates with lower risk of GL (14). Higher variability in pulse pressure correlated with higher odds of 1- and 3-year GL, and higher pulse rates correlated with 3-year GL and mortality in the first 3 years, likely reflecting nonresolution of the microvascular milieu in the renal vasculature and heightened sympathetic drive underlying chronic kidney disease and its associated cardiovascular risk (15–18). Transplants with better eGFRs and a continued upward trend in eGFR in the first 90 days posttransplant as well as those with a continued upward trend in hemoglobin were significantly associated with lower risk of GL. We were thus able to capture evolution of allograft function in a granular manner by incorporating eGFR as a trajectory variable in our models. Similar findings have been reported by other groups using traditional modeling approaches (19). These findings relating to eGFR and resolution of anemia of renal disease further reflect a robust resolution of uremic pathobiology as a harbinger of better graft survival. Thus, our approach brings an auxometric (variables reflecting growth or change in biology) dimension to modeling
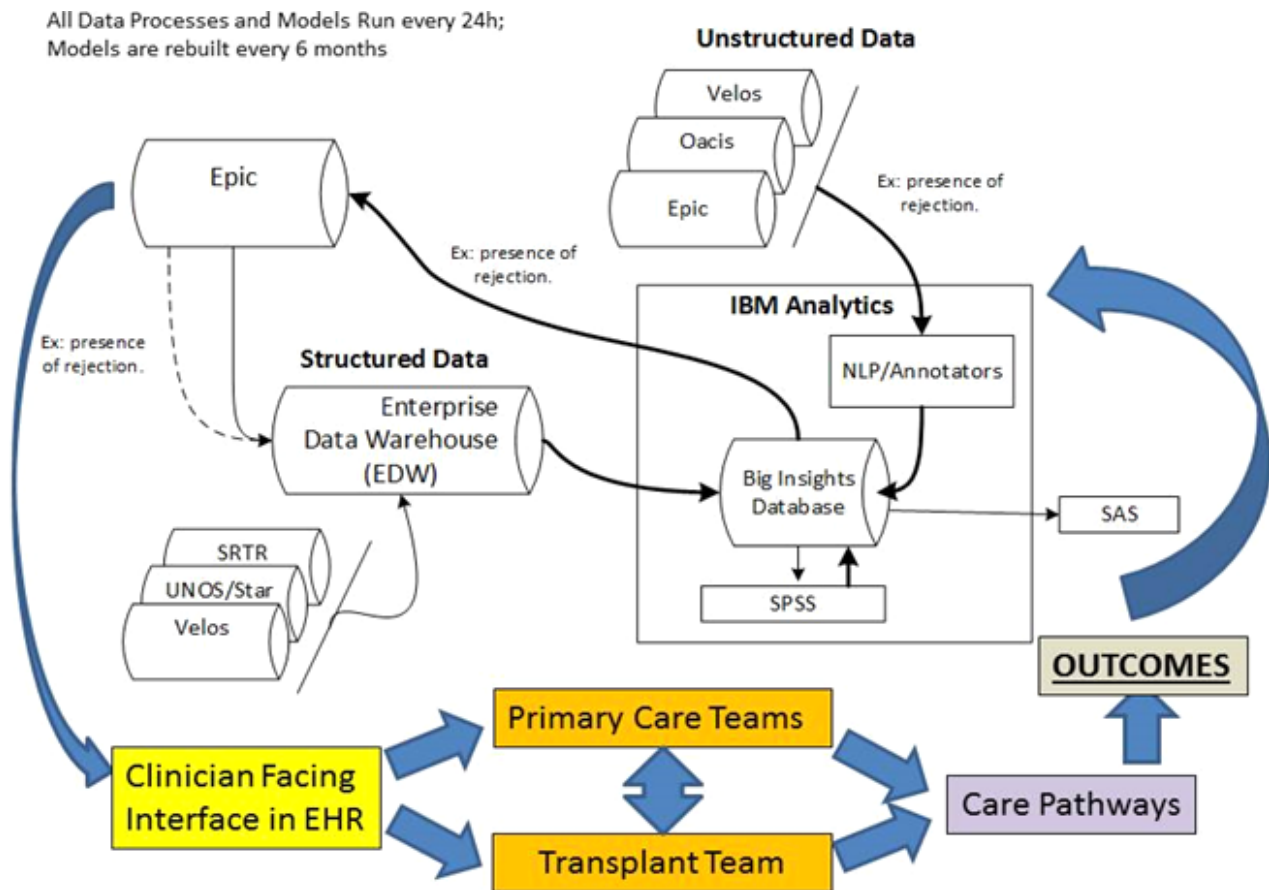
**Figure 2: Comparison of four predictive models for graft loss risk among kidney transplants.**\* \*Model 1: UNOS data only. Model 2: UNOS + Tx database. Model 3: UNOS + Tx database + EHR comorbidity. Model 4: UNOS + Tx database + EHR comorbidity + EHR posttransplant trajectory + NLP. AUC, area under the curve; EHR, electronic health record; NLP, natural language processing; Tx, transplant; ROC, receiver operating characteristic curve; UNOS, United Network for Organ Sharing.

approaches as was first articulated by Feinstein and colleagues in 1974: The addition of functional substrates that reflect biology adds to predictive accuracy (8).

The association between evolving immunologic risk and GL was captured within our models. The accompaniment of underimmunosuppression, acute rejection, was associated with increased risk of 3-year GL. In this regard, the relationship between viral infections, a correlate of

overimmunosuppression in these models, was intriguing. CMV infection was protective for early GL, which may suggest that clinicians are reacting to these early CMV infections by significantly reducing immunosuppression and likely reducing the risk of early GL (20).

Despite the primacy of rejection prophylaxis in the medical management of the transplant recipient and the relevance of histologic lesions to outcome, large national databases

**Figure 3: Data process flow**. Schema for deployment of the predictive model for graft loss in the clinic is depicted. Transplant and Primary care teams can synergize using a data-driven daily deployable EHR embedded interface that curates multiple data sources and drives clinical workflow using predictive modeling. Our solution allows transplant programs to optimize outcomes proactively by putting data to drive workflow and allocate resources rather than reacting to regulatory stress (Oacis: Electronic Health Record Pre-2012; Epic: Electronic Medical Record; Velos: Proprietary Transplant Database; EDW: Enterprise Data Warehouse). This workflow can be customized to transplant centers' unique patient populations and can be updated with new releases of SRTR models as well. EHR, electronic health record; SRTR, Scientific Registry of Transplant Recipients.

are remarkably limited in their capture of rejection, its treatment, and postrejection evolution of renal function (2,21). We were able to surmount this deficiency by incorporating data on histologic lesion scores from pathologic reports on allograft biopsies using NLP techniques.

Cardiovascular disease is a major driver of posttransplant mortality and costs. Our approach identified this association and its contribution to the force of mortality in a direction consistent with prior experiences. Additionally, the results underscore the importance of social support and a greater degree of educational attainment as being protective for GL, which is consistent with our prior observations (22,23).

Taken together, our data suggest that several windows of opportunity exist for action on mutable domains such

as cardiovascular disease, immunologic risk, or social determinants of health. These could be better identified earlier through prospective application of our predictive solution in the clinical workflow. The event-dense and data-rich clinical environment of transplantation is served by a manpower structure that is already well positioned with care coordination resources. Such an approach can be customized to the unique demographic, clinical, and socioeconomic characteristics of patients at individual transplant centers so as to accurately identify groups of patients at risk for GL or death.

The overall concept of such a process flow as applicable to kidney transplant recipients is presented in Figure 3 and accompanying legend. We are currently building EHR-based workflows where data structures and analytic solutions power an automated daily capture of model

variables, calculate and update individual patient risk (probability of event), and push these model-derived risk scores in a clinician-facing interface that is quotidian refreshed and drives clinical practice. This workflow is aimed at triggering actions by appropriately skilled teams as opposed to making all tasks part of the transplant team's workflow. Such an approach has been used successfully to reduce readmissions among heart failure patients (24).

There are limitations of our modeling approach that require discussion. The study is limited to adults (age ≥18), and this population may not represent the characteristics of patients seen in other community-based or academic-based transplant centers. This analysis describes the kidney transplant population and risk assessment in a single academic-affiliated hospital; external validity of this model or included variables has not been conducted. In addition, utilization may be underestimated, since patients can use any emergency department or hospital in the region. Matching our data to a larger universal billing claims database or prescription databases could minimize missing utilization (25), as we defined utilization by event numbers, not resource consumption, care cost, or charges. It should be noted that the SRTR uses survival methodology in its modeling of outcomes. In this regard, results of our logistic modeling approach are directionally consistent with survival approaches applied to our data. We further submit that our major aim was to provide proof of concept of a pragmatic approach to the use of data available in the EHR to guide optimal and timely clinical interventions and not necessarily an attempt to uncover novel associations beyond those discernible with national data. Furthermore, our results should not be regarded as presenting the ideal or definitive prediction model but more as a method by which centers could customize their data structures to serve in clinically actionable ways. As such, external validation of our approach is needed and welcome.

## Conclusions and Future Directions

Inclusion of dynamically evolving structured and unstructured longitudinal patient-level data using Big Data approaches improves the accuracy of prediction of GL and mortality among renal transplant recipients. Through electronic capture and curation, this approach to predictive modeling can be feasibly automated and deployed to provide near real-time clinically actionable results that have the potential to optimize graft and patient outcomes. Such a data structure driving predictive analytics to the bedside has the potential to empower clinicians to optimize value through precision care delivery in transplantation. The unique data-driven care models so developed in transplantation could impact clinically relevant high-value use cases in fields of medicine outside transplantation.

## Disclosure

The authors of this manuscript have no conflicts of interest to disclose as described by the *American Journal of Transplantation*.

## References

1. Meier-Kriesche HU, Schold JD, Srinivas TR, Kaplan B. Lack of improvement in renal allograft survival despite a marked decrease in acute rejection rates over the most recent era. Am J Transplant 2004; 4: 378–383.
2. Kaplan B, Schold J, Meier-Kriesche HU. Overview of large database analysis in renal transplantation. Am J Transplant 2003; 3: 1052–1056.
3. Taber DJ, Palanisamy AP, Srinivas TR, et al. Inclusion of dynamic clinical data improves the predictive performance of a 30-day readmission risk model in kidney transplantation. Transplantation 2015; 99: 324–330.
4. McAdams-Demarco MA, Grams ME, King E, Desai NM, Segev DL. Sequelae of early hospital readmission after kidney transplantation. Am J Transplant 2014; 14: 397–403.
5. IBM White Paper. 5 Steps to becoming a data-driven healthcare organization. Somers, NY: IBM Corporation, 2016; p. 1–7.
6. Racusen LC, Solez K, Colvin RB, et al. The Banff 97 working classification of renal allograft pathology. Kidney Int 1999; 55: 713–723.
7. Quan H, Sundararajan V, Halfon P, et al. Coding algorithms for defining comorbidities in ICD-9-CM and ICD-10 administrative data. Med Care 2005; 43: 1130–1139.
8. Charlson ME, Feinstein AR. The auxometric dimension. A new method for using rate of growth in prognostic staging of breast cancer. JAMA 1974; 228: 180–185.
9. Firth D. Bias reduction of maximum likelihood estimates. Biometrika 1993; 80: 27–38.
10. Heinze GA. A solution to the problem of separation in logistic regression. Stat Med 2002; 21: 2409–2419.
11. Heinze GA. A comparative investigation of methods for logistic regression with separated or nearly separated data. Stat Med 2006; 25: 4216–4226.
12. Harrell FE Jr, Lee KL, Mark DB. Multivariable prognostic models: Issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. Stat Med 1996; 15: 361–387.
13. SRTR risk adjustment model documentation: Waiting list and post-transplant outcomes. 2016 [cited 2016 July16]. Available from: http://www.srtr.org/csr/current/modtabs.aspx.
14. Rao PS, Schaubel DE, Guidinger MK, et al. A comprehensive risk quantification score for deceased donor kidneys: The kidney donor risk index. Transplantation 2009; 88: 231–236.
15. Amann K, Wanner C, Ritz E. Cross-talk between the kidney and the cardiovascular system. J Am Soc Nephrol 2006; 17: 2112–2119.
16. Chang TI, Tabada GH, Yang J, Tan TC, Go AS. Visit-to-visit variability of blood pressure and death, end-stage renal disease, and cardiovascular events in patients with chronic kidney disease. J Hypertens 2016; 34: 244–252.

17. Johnson RJ, Rodriguez-Iturbe B, Kang DH, Feig DI, Herrera-Acosta J. A unifying pathway for essential hypertension. Am J Hypertens 2005; 18: 431–440.

18. Meier-Kriesche HU, Schold JD, Srinivas TR, Reed A, Kaplan B. Kidney transplantation halts cardiovascular disease progression in patients with end-stage renal disease. Am J Transplant 2004; 4: 1662–1668.

19. Wan SS, Cantarovich M, Mucsi I, Baran D, Paraskevas S, Tchervenkov J. Early renal function recovery and long-term graft survival in kidney transplantation. Transpl Int 2016; 29: 619–626.

20. Elfadawy N, Flechner SM, Liu X, et al. CMV Viremia is associated with a decreased incidence of BKV reactivation after kidney and kidney-pancreas transplantation. Transplantation 2013; 96: 1097–1103.

21. Gonzales MM, Bentall A, Kremers WK, Stegall MD, Borrows R. Predicting individual renal allograft outcomes using risk models with 1-year surveillance biopsy and alloantibody data. J Am Soc Nephrol 2016. [Epub ahead of print].

22. Goldfarb-Rumyantzev AS, Rout P, Sandhu GS, Khattak M, Tang H, Barenbaum A. Association between social adaptability index and survival of patients with chronic kidney disease. Nephrol Dial Transplant 2010; 25: 3672–3681.

23. Taber DJ, Hamedi M, Rodrigue JR, et al. Quantifying the race stratified impact of socioeconomics on graft outcomes in kidney transplant recipients. Transplantation 2015; 100: 1550–1557.

24. Evans RS, Benuzillo J, Horne BD, et al. Automated identification and predictive tools to help identify high-risk heart failure patients: Pilot evaluation. J Am Med Inform Assoc 2016; 23: 872–878.

25. Massie AB, Kucirka LM, Segev DL. Big data in organ transplantation: Registries and administrative claims. Am J Transplant 2014; 14: 1723–1730.

## Supporting Information

Additional Supporting Information may be found in the online version of this article.

**Table S1:** One-year graft loss.

**Table S2:** Three-year graft loss.

**Table S3:** Three-year mortality.

**Data S1:** Explanatory notes on methods.