## 1 A New Approach to Develop Large Scale Land Use Models Using Open Source Data

2 3

#### 3 4

## 5 Ali Riahi Samani

- 6 Department of Civil Engineering
- 7 University of Memphis, Memphis, Tennessee, United States, 38152
- 8 Email: <u>rhsamani@memphis.edu</u>
- 9

# 10 Sabyachee Mishra

- 11 Department of Civil Engineering
- 12 University of Memphis, Memphis, Tennessee, United States, 38152
- 13 Email: smishra3@memphis.edu
- 14

## 15 David Jung-Hwi Lee

- 16 Assistant Director, Long Range Planning
- 17 Tennessee Department of Transportation
- 18 Nashville, Tennessee, 37243-0344
- 19 Email: <u>david.Lee@tn.gov</u>
- 20

## 21 Mihalis M. Golias

- 22 Department of Civil Engineering
- 23 University of Memphis, Memphis, Tennessee, United States, 38152
- 24 Email: <u>mgkolias@memphis.edu</u>
- 25

## 26 Jerry Everett

- 27 Center for Transportation Research
- 28 The University of Tennessee Knoxville, Tennessee, United State, 37996-4133
- 29 Email: jeverett@utk.edu

# 1 ABSTRACT

2 Developing a land-use model for large-scale cases is a topic that has received less attention in the literature 3 while transportation engineers and urban planners continue to analyze the effect of various policies in multi-4 jurisdiction metropolitan areas and to some extent in statewide scale. Gravity based models when too 5 simplistic, microsimulation models require extensive data and massive computation. This paper presents a 6 land use model that can be applied to large-scale geographies using open source data and be able to forecast 7 demographic and socioeconomic attributes with reasonable accuracy and acceptable computational time. 8 The proposed model incorporates the Putman's Integrated Transportation-Land Use Package (TELUM) 9 and Kockelman's Gravity-based Land Use Model (G-LUM) fundamentals with enhanced formulation of 10 newly added variables and structural changes. Considering the non-convex and non-linear nature of the 11 proposed model, we utilize an enhanced genetic algorithm for base year calibration. Further, we assess the accuracy of the model with two-fold validation including back-casting and forecasting. We utilize the state 12 13 of Tennessee as the case study area and utilized all open source data available to the model application. The 14 model results show reasonably accurate estimates of households by size, employment by industry, and land 15 utilization by condition. As applicable the model outperforms G-LUM and TELUM by accuracy ( $\mathbb{R}^2$ ) and 16 error measures (MAPE). The proposed land use model has the potential to be applied for medium to large 17 scale geographies with reasonable accuracy in predicting socio-economic, demographic, and land condition

18 estimates by using open source data.

19 Keywords: Integrated Land Use, Transport Model, Gravity Theory, Statewide Land Use Model, Genetic

20 Algorithm

#### 1 INTRODUCTION

2 The interdependence of land use and transportation led urban planners and transportation 3 engineering to focus on each while the former has received less attention in recent years. The first generation 4 of land use models such as aggregate spatial interaction and gravity models were introduced around 1960s. 5 Lowry was a pioneer by introducing the model of metropolis (1). Later on, utility-based econometric and 6 discrete choice models were developed. The development of advanced micro-simulation land use models 7 and activity-based travel demand models created the need for a new generation of integrated land use-8 transport systems (2). New models such as ILUTE (3) and ILUMASS (4) were developed thereafter. 9 Existing models such as UrbanSim (5), PECAS (6), and MUSSA (7) were updated to facilitate the need for 10 advanced research in the field of integrated land use-transport modeling.

11 In recent years, improving the accuracy of the land use models by introducing micro-simulation 12 garnered much attention because of behavioral interpretation, studies on their accuracy is still an evolving 13 area of research. The problem of dealing with large-scale geographies received less attention though 14 multiple jurisdictions are managed by metropolitan planning authorities, or state planning agencies. T the 15 topic of accuracy versus scale became important on the regional scale, such as multi-jurisdiction 16 metropolitan areas and in statewide applications. Despite the high accuracy of microsimulation land use 17 models (i.e. UrbanSim (8)), the enormous data requirement and massive computation time, make the 18 implementation of these models on large-scale cases challenging (9). The development of a large-scale 19 land use model becomes more important when the integration of land use models with a statewide travel 20 demand model is raised. It is important to urban planners and transportation engineers to assess and analyze 21 the effect of policies or scenarios on broader scales. Therefore, the development of a land use model can be 22 applied in large-scale (regional or statewide) becomes crucial.

23 The purpose of this research is to develop a land use model which can be applied to large-scale 24 geographies with acceptable computational time reasonable forecasting accuracy and use of open-source 25 data. Such large-scale models can be integrated with existing travel demand models as applicable in many areas. Few land use models are applied in large-scale geographies, such as PECAS and TELUM (10). 26 27 PECAS is a generalized approach for simulating spatial economic systems. It operates by clearing spatial 28 submarkets for various goods, services, and factors in a short-run equilibrium based on development event 29 probabilities (6). California's statewide land use model is a statewide PECAS model, integrated with a 30 statewide travel demand forecast model (11). Moreover, PECAS components applied in the development 31 of statewide transportation land use modeling systems for Ohio and Oregon (12). However, the 32 implementation of PECAS to large-scale cases has limitations on computational time and the number of 33 zones (13).

34 TELUM is another well-known land use model. TELUM is an integrated land use and transport 35 model that incorporates gravity theory to allocate households and employment to zones (14). TELUM 36 development is based on three components, a disaggregated residential allocation model (DRAM), an 37 employment allocation model (EMPAL), and a land consumption model (LANCON) (10). Presenting a 38 user-friendly interface, GIS-base result, reasonable data requirements, and short run time have gained 39 researchers' attention to applying this model in their projects (15–19). Although the implementation of 40 TELUM is simple, this model has limitations. First, in TELUM the number of employment categories is 41 limited to 4 to 8 employment sections and the number of household categories is limited to 5-8 categories, usually for household's income. Second, TELUM has restrictions on zone size and it is recommended that 42 43 the average population in a zone lies between 3,000 and 10,000 (10).

44 Kockelman et al. tried to solve these limitations of TELUM by developing Gravity Land Use 45 Model (G-LUM) (20). The G-LUM structure is based on the formulation of the ITLUP package (21) and includes three major sub-models for predicting changes in employment location (EMPLOC), residential 46 47 location (RESLOC), and land consumption (LUDENSITY). G-LUM was used to validate the outputs of 48 TELUM. Studies showed that, although G-LUM and TELUM use the same structure, the forecasting results 49 were varied (22). The differences rooted in the method of calibration (19). TELUM uses a gradient search 50 method and G-LUM uses the Nelder-Mead method with 12 different initial points. In addition to the 51 mentioned limits, the formulation of land use consumption in TELUM and G-LUM (LANCON and

1 LUDENSITY) tends to generate unreasonable average land consumption values for households and jobs 2 (as compared with base and prior year land conditions in each zone) (23).

The contribution of this paper is three-fold. First, improvement in the model structure of TELUM and G-LUM, so that large scale implementation is possible. Second, proposition of a new solution algorithm to enhance the accuracy of the land use models. Third, demonstration of model applicability using a case study area by only using open source data. The rest of the paper is organized as follows. The next section discusses the proposed methodology and calibration procedure. The following section presents the data used and the case study. The results section compared the performance of the proposed models with similar models in the past. The conclusion section summary of the paper, and avenues for future research.

10

# 11 METHODS

In this section, the proposed land use model development and specifications are discussed. First, the description and formulation of different sections of the model are provided. Then, a brief description of the travel demand model which is planning to be integrated with the proposed land use model is provided. Finally, the calibration of the proposed model is discussed. As a general description, the proposed model is based on gravity theory and incorporates the principle structure of the TELUM and G-LUM, while the details have been changed to improve the model's accuracy. The model presented in TAZ level and forecasts

18 socioeconomic and demographic character of zones in five-year intervals.

## 19 Land Use Models

In this section, two land use models' structure are provided. The first model is called Large-Scale Land Use Model (LS-LUM) and the second model is named Large-Scale Land Use Model Without House Condition (LS-LUM-WOHC). In the following subsections, the formulation of both models is provided.

23 Large-Scale Land Use Model (LS-LUM)

24 LS-LUM contains two principal sections and two subsections. Principal sections estimate 25 households and employments in different categories. These two models incorporate gravity theory to 26 allocate households and employments to each TAZ. The first principal model is named HH-AL (Households 27 Allocation) which is responsible for residential location choice. This model assigns households to each 28 TAZ based on the total number of houses, vacant houses, the amount of residential land (acres), total useable 29 land in each TAZ, and the travel cost between zones. The second principal model is called EMP-AL 30 (Employments Allocation). This model allocates employments to zones based on job opportunities in the prior year (lag year), the amount of commercial, industrial, and agricultural land (acres) in each TAZ, and 31 32 the travel cost between zones. Two subsections are responsible for updating house conditions and land use 33 consumption which are the components of principal sections. Two models provided in these subsections 34 are called HC (House Condition) and LC (Land Consumption). HC models the number of total and vacant 35 houses in each TAZ and LC Models the amount of land (acres) in each land use class (residential, 36 commercial, industrial, agricultural, and developable or vacant). Multiple Linear Regression (MLR) is 37 applied in these subsections. In the following sections, the formulations of these models are provided. In general, in this paper i and j represent TAZs,  $C_{i,i}$  represents the travel cost between zone i and zone j, and 38 t represents the period of time (i.e. year 2010). Moreover, n and k stand respectively for household 39 40 categories (i.e. different household size) and employments categories (i.e. NAICS sectors categories).

42

$$N_{i,t}^{n} = \eta^{n} \sum_{j} a_{n} E_{j,t}^{T} \frac{W_{i,t-1}^{n} C_{i,j,t-1}^{\alpha^{n}} \exp\left(\beta^{n} C_{i,j,t-1}\right)}{\sum_{i} W_{i,t-1}^{n} C_{i,j,t-1}^{\alpha^{n}} \exp\left(\beta^{n} C_{i,j,t-1}\right)} + (1 - \eta^{n}) N_{i,t-1}^{n}$$
(1)

43

44

Where

$$W_{i,t-1}^{n} = (H_{i,t-1}^{T})^{o^{n}} (H_{i,t-1}^{V})^{p^{n}} (1 + \frac{L_{i,t-1}^{Res}}{L_{i,t-1}^{T}})^{q^{n}}$$
(2)

In Equation 1,  $N_{i,t}^n$  is the number of households in category *n* in zone *i* in time *t*,  $a_n$  is the proportion of the population to employment in zone *i*,  $E_{j,t}^T$  is the total number of employments,  $W_{i,t-1}^n$  is the attractiveness function of zone *i* to which attract employment in zone *j* to live in zone *i* in year t - 1.  $W_{i,t-1}^n$  is a weighted multiplication of different components in a zone. In Equation 2,  $H_{i,t-1}^T$  is the total number of houses,  $H_{i,t-1}^v$  is the number of vacant houses,  $LP_{i,t-1}$  is the residential land value,  $L_{i,t-1}^{Res}$  is the amount of residential land are in zone *i* in year t - 1. Finally,  $\eta$ ,  $\alpha$ ,  $\beta$ , o, p, and q are parameters estimated in the calibration procedure.

$$E_{j,t}^{k} = \lambda^{k} \sum_{i} N_{i,t-1}^{T} \frac{M_{j,t-1}^{k} C_{i,j,t-1}^{\alpha^{k}} \exp\left(\beta^{k} C_{i,j,t-1}\right)}{\sum_{i} M_{j,t-1}^{n} C_{i,j,t-1}^{\alpha^{k}} \exp\left(\beta^{k} C_{i,j,t-1}\right)} + (1 - \lambda^{n}) E_{j,t-1}^{k}$$
(3)

10

11 Where,

$$M_{j,t-1}^{k} = (E_{j,t-1}^{k})^{g^{k}} (L_{j,t-1}^{Com} + L_{j,t-1}^{Ind} + L_{j,t-1}^{Agr})^{h^{k}}$$
(4)

12

13 In **Equation 3**,  $E_{j,t}^{k}$  is the number of employments in category k,  $N_{i,t-1}^{T}$  is the total number of 14 households, and  $M_{j,t-1}^{k}$  is the attractiveness function shows how much zone j is attractive for peopled living 15 in zone i to find a job.  $M_{j,t-1}^{k}$  is calculated based on, job opportunities in year t - 1 ( $E_{j,t-1}^{k}$ ), amount of 16 commercial ( $L_{j,t-1}^{com}$ ), industrial ( $L_{j,t-1}^{Ind}$ ), and agricultural ( $L_{j,t-1}^{Agr}$ ) land in zone j in year t - 1.  $\lambda$ ,  $\alpha$ ,  $\beta$ , g, and 17 h are parameters estimated in the calibration procedure.

In this subsection, the total number of houses and the number of vacant houses in each TAZ are updated. First, the total number of houses in each TAZ is calculated by applying a Multiple Linear Regression. As **Equation 5** shows, the total number of houses in zone *i* and in year  $t(H_{i,t}^T)$  is the dependent variable; while, the number of total houses in the previous year (t - 1), the amount of vacant land  $(L_{i,t-1}^{Vac})$ , and the total number of households are the independent variables.

$$H_{i,t}^{T} = \theta_{0} + \theta_{1}(H_{i,t-1}^{T}) + \theta_{2}(H_{i,t-1}^{V}) + \theta_{3}(L_{i,t-1}^{Vac}) + \theta_{4}(N_{i,t}^{T}) + \varepsilon$$
(5)

In Equation 5,  $L_{i,t-1}^{Vac}$  is the amount of vacant or developable land in zone *i* and  $\varepsilon$  is the error associated in regression. In this equation,  $\theta_0$  is the intercept and  $\theta_1$  to  $\theta_4$  are coefficient estimated in calibration.

After calculating the total number of houses, the number of vacant houses in each TAZ can be estimated as follow:

$$H_{i,t}^{V} = H_{i,t}^{T} - \sum_{n} N_{i,t}^{n}$$
(6)

29

30 *4. LC*:

Finally, in *LC*, the amount of land in different land use classes is updated to feed the two principal models (*HH-AL* and *EMP-AL*) in order to forecast future years' demographic and socio-economic conditions.

$$L_{i,t}^{Res} = R_0 + R_1 L_{i,t-1}^{Vac} + R_2 (L_{i,t-1}^{Res}) + R_3 (N_{i,t-1}^T) + R_4 (N_{i,t}^T) + \varepsilon$$
(7)

$$L_{i,t}^{Com} = C_0 + C_2(L_{i,t-1}^{Vac}) + C_2(L_{i,t-1}^{Com}) + C_3(E_{i,t-1}^{Com}) + C_4(E_{i,t}^{Com}) + \varepsilon$$
(8)

$$L_{i,t}^{Ind} = I_0 + I_1(L_{i,t-1}^{Vac}) + I_2(L_{i,t-1}^{Ind}) + I_3(E_{i,t-1}^{Ind}) + I_4(E_{i,t}^{Ind}) + \varepsilon$$
(9)

$$L_{i,t}^{Agr} = A_0 + A_1 \left( L_{i,t-1}^{Vac} \right) + A_2 \left( L_{i,t-1}^{Agr} \right) + A_3 \left( E_{i,t-1}^{Agr} \right) + A_4 \left( E_{i,t}^{Agr} \right) + \varepsilon$$
(10)

$$L_{i,t}^{Vac} = L_{i,t-1}^{Vac} - \left(L_{i,t-1}^{Res} - L_{i,t}^{Res}\right) - \left(L_{i,t-1}^{Com} - L_{i,t}^{Com}\right) - \left(L_{i,t-1}^{Ind} - L_{i,t}^{Ind}\right) - \left(L_{i,t-1}^{Agr} - L_{i,t}^{Agr}\right)$$
(11)

In **Equations 9** to **10**,  $E^{Arg}$  refers to the number of employments in NAICS sector 11 (agriculture, forestry, fishing, and hunting),  $E^{Com}$  is the number of employments in NAICS sectors 44, 45, 51, 52, 53, and 72 (retail trade, finance and insurance, real estate and rental and leasing, accommodation and food services), and  $E^{Ind}$  is the number of employments in NAICS sectors 21, 31, 33, and 42 (mining, quarrying, oil and gas extraction, manufacturing, and wholesale trade).

#### 9 Large-Scale Land Use Model Without House Condition (LS-LUM-WOHC)

10 LS-LUM-WOHC is the second model developed in this paper. The formulation of this model is 11 very similar to LS-LUM and the only difference is that in this land use mode, the house condition subsection 12 (HC) and its components have been removed from LS-LUM. The purpose of developing this land use model 13 is to evaluate the effect of adding HC to land use modeling. In other words, developing LS-LUM-WOHC 14 provides the opportunity to compare the presence and absence of HC. Therefore, LS-LUM-WOHC is consist of two principal models and a subsection model. The principal models are responsible for allocating 15 16 households and employment. LS-LUM-WOHC incorporates the same model for allocating employment. 17 *EMP-Al* is applied in this model too. However, the household allocation model is different in comparison 18 with LS-LUM. Household allocation model in LS-LUM-WOHC is called HH-AL2, where the formulation 19 is as follow:

$$N_{i,t}^{n} = \eta^{n} \sum_{j} a_{n} E_{j,t}^{T} \frac{W_{i,t-1}^{n} C_{i,j,t-1}^{a^{n}} \exp\left(\beta^{n} C_{i,j,t-1}\right)}{\sum_{i} W_{i,t-1}^{n} C_{i,j,t-1}^{a^{n}} \exp\left(\beta^{n} C_{i,j,t-1}\right)} + (1 - \eta^{n}) N_{i,t-1}^{n}$$
(12)

20 Where,

$$W_{i,t-1}^{\prime n} = \left(1 + \frac{L_{i,t-1}^{Res}}{L_{i,t-1}^{T}}\right)^{q^{\prime n}}$$
(13)

In the formulation of *HH-AL2*, the attractiveness of each zone  $(W'_{i,t-1}^n)$  is calculated only by using the

amount of residential  $(L_{i,t-1}^{Res})$  and the total land  $(L_{i,t-1}^{T})$  in zone *i* and in the prior year. In addition, the

formulation for the subsection model is similar to LS-LUM; where, *LC* is applied to forecast the amount of

24 land in five different land use classes.

#### 1 Model Explanation

2 This integrated modeling framework starts with forecasting employment in different categories and for each 3 zone (see Figure 1). This section of the model gets employment, the amount of agricultural, commercial, 4 industrial lands, and travel cost in each zone and for the prior year. The output of this section is the 5 forecasted employment (by different categories) in each TAZ. The output of the EMP-AL would serve as 6 input for the HH-AL. The HH-AL incorporates the current total employment (from EMP-AL), the total 7 number of houses, the number of vacant houses, and the proportion of residential to total land in each zone 8 for the prior year. The output of this section is the number of households (by different categories, e.g. 9 income). Then HC computation is processed, by forecasting how many houses will be built in each TAZ. 10 This section needs total and vacant number of houses, the amount of vacant land in the prior year, and the forecasted total number of households (from HH-AL section). Considering the total number of forecasted 11 12 houses and total households in each TAZ, HC forecasts the number of vacant houses in each zone by 13 subtracting the total number of households from the total number of houses in each zone. The output of HC14 feeds the HH-AL by providing the number of total and vacant houses. Lastly, LC forecasts the amount of 15 residential, commercial, industrial, agricultural, and vacant land (developable) in each zone and for each forecasting year. The output of LC directly affects other models' sections. By connecting LC to other 16 17 sections, capturing the effect of land use changes on the socio-economic character of each TAZ would be 18 possible and more accurate results can be obtained. The amount of commercial, industrial, and agricultural 19 land modeled in this section are added to employment section. The amount of residential land is added to 20 the HH-AL. Finally, the amount of vacant land is one of the components involved in forecasting the total 21 number of houses in a zone. Moreover, the amount of vacant land in each zone works as a development 22 restriction. Because in the model, if all the vacant land had allocated to other land use classes, no more 23 development will happen, and the model will stop adding a new area to other land use classes (residential, 24 commercial, industrial, and agricultural).



#### Figure 1 Integrated land use transport model's flowchart (dashed lines represent one period (t - 1)lagged feedback of information; each period is 5 years).

#### **3 Travel Demand Model**

4 The Travel Demand Model (TDM), which integrated with the land use model and the travel time 5 derived form, is the Tennessee Statewide Travel Model (TSTM) version 3 (24). This version of TSTM is a 6 traditional four-step, TDM consisting of three different components, short distance passenger model (trips 7 less than 50 miles), long-distance passenger model, and freight model. The underlying geographic area of 8 operation is at the TAZ level. The total number of TAZs in TSTM is 3,687. Zonal attributes include the 9 number of households, categorized by income, size, worker, presence of student, presence of seniors, and 10 the number of vehicles; and the number of employments categorized by 20 sectors of NAICS codes. The TSTM3 can be understood at a high level as comprised of input network and socioeconomic data together 11 12 with some component demand models and a highway assignment model. The demand components can be 13 gathered in three broad groups related to short-distance passenger demand, long-distance passenger 14 demand, and freight and truck demand. The TSTM3 uses TransCAD's implementation of the tri-conjugate 15 Frank-Wolfe algorithm for multi-class user equilibrium traffic assignment (25). The accessibility matrices 16 which serve as input for the land use model are obtained from TSTM's assigned networks using shortest 17 path method.

18

## 19 Calibration

The parameters of four models (*HH-AL*, *EMP-AL*, *HC*, and *LC*) need to be estimated through a calibration process. The calibration of the proposed models is categorized into two sections. The first section is dedicated to the estimation of the parameters of *HC* and *LC*. These two models are Multiple Linear Regression and the intercept and coefficients are estimated using least square method. The objective of the second section is to estimate the parameters of *HH-AL* and *EMP-AL*. The calibration is conducted through maximum likelihood approach where the two following objective functions are defined. First, for *HH-AL* the objective function is as below:

$$Z_{1} = Min \sum_{i} \sum_{n} \frac{\left(N_{i,t_{Obs}}^{n} - N_{i,t_{Est}}^{n}\right)^{2}}{\left(\sigma_{P_{i,t_{Obs}}^{n}}\right)^{2}}$$

(14)

27 Where, 
$$N_{i,t_{Est}}^{n}$$
 is defined in equation (1), (2), and is illustrated here for convenience.

$$N_{i,t_{Est}}^{n} = \eta^{n} \sum_{j} a_{n} E_{j,t}^{T} \frac{W_{i,t-1}^{n} C_{i,j,t-1}^{\alpha^{n}} \exp\left(\beta^{n} C_{i,j,t-1}\right)}{\sum_{i} W_{i,t-1}^{n} C_{i,j,t-1}^{\alpha^{n}} \exp\left(\beta^{n} C_{i,j,t-1}\right)} + (1 - \eta^{n}) N_{i,t-1}^{n}$$
(1)

28

$$W_{i,t-1}^{n} = (H_{i,t-1}^{T})^{o^{n}} (H_{i,t-1}^{V})^{p^{n}} (1 + \frac{L_{i,t-1}^{Res}}{L_{i,t-1}^{T}})^{q^{n}}$$
(2)

29

30 In the objective function  $Z_1$ ,  $N_{i,t_{Obs}}^n$  and  $N_{i,t_{Est}}^n$  are respectively, the number of observed and 31 estimated households in category *n* and zone *i* and  $\sigma_{P_{i,t_{Obs}}^n}$  is the standard deviation of observations. Where, 32 the decision variables are  $\eta$ ,  $\alpha$ ,  $\beta$ , o, p, and q (the calibration parameter mentioned in **Equations 1** and **2**). 33 Moreover, a similar objective function is defined for *EMP-AL* as follows:

$$Z_{2} = Min \ \sum_{j} \sum_{k} \frac{\left(E_{j,t_{Obs}}^{k} - E_{j,t_{Est}}^{k}\right)^{2}}{\left(\sigma_{E_{j,t_{Obs}}^{k}}\right)^{2}}$$
(15)

34

Where,  $E_{j,t_{ESt}}^k$  is defined in equation (1), (2), and is illustrated here for convenience.

$$E_{j,t}^{k} = \lambda^{k} \sum_{i} N_{i,t-1}^{T} \frac{M_{j,t-1}^{k} C_{i,j,t-1}^{\alpha^{k}} \exp\left(\beta^{k} C_{i,j,t-1}\right)}{\sum_{i} M_{j,t-1}^{n} C_{i,j,t-1}^{\alpha^{k}} \exp\left(\beta^{k} C_{i,j,t-1}\right)} + (1 - \lambda^{n}) E_{j,t-1}^{k}$$
(3)

2

$$M_{j,t-1}^{k} = (E_{j,t-1}^{k})^{g^{k}} (L_{j,t-1}^{Com} + L_{j,t-1}^{Ind} + L_{j,t-1}^{Agr})^{h^{k}}$$
(4)

Similarly, in objective function  $Z_2$ ,  $E_{j,t_{Obs}}^k$  and  $E_{j,t_{Est}}^k$  are, the number of observed and estimated employments in category k and zone j and  $\sigma_{E_{j,t_{Obs}}^k}$  is the standard deviation of observations respectively.

## 5 Where, the decision variables are $\alpha$ , $\beta$ , g, and h, defined in Equations 3 and 4.

6 Both objective functions  $Z_1$  and  $Z_2$  are non-linear, non-convex, and are not subjected to any 7 constraints. In the previous land use models (TELUM and G-LUM), a gradient search method and the 8 Nelder-Mead method with 12 different initial points applied. Previous approaches add strict limitations to 9 the solution approach. First, due to the non-convexity of objective functions, using gradient search method 10 would increase the chance of trapping in local optimum solution. Second, accuracy and final solution of the 11 Nelder-Mead method with initial points is highly sensitive to selection of initial points (23). Therefore, in 12 order to eliminate these limitations, in this paper, an evolutionary algorithm is applied to solve the above-13 mentioned optimization problem. The following section discusses the proposed solution approach.

14

#### 15 Solution Approach

- 16 In this paper, a Modified Genetic Algorithm (GA) is applied to solve these optimization problems. GA
- 17 applies to solve non-convex and non-linear optimization problems because of its superiority in evolutionary
- 18 search computation over other search techniques which are limited by the continuity, differentiability, and
- unimodality of the evaluated functions (26). GA operates by maintaining and modifying the characteristics
   of a set of trial solutions (population) over iterations (generations). Each of the GA steps is further illustrated
   below:
- Encoding: the initial step in operating the genetic algorithms is forming an initial population (initial trial solution set). Each individual solution in the population is represented by a binary string which is called chromosome. Each chromosome contains model parameters that are encoded in the form of binary codes (called genes). In the initial set, the values of the model parameters are randomly assigned. In this research, the model generates 1,000 chromosomes in the initial population.
- Reproduction: the initial population will not provide an optimal solution. A genetic algorithm works by trying to reproduce other chromosomes that are better solutions. The reproduction process is simply a selection process where those chromosomes that have a better objective function will have a higher chance of reproduction. Through this process, the overall quality of the population will gradually be improved.
- Crossover: in the crossover, genetic materials (genes) between chromosomes are exchanged to generate a new chromosome. Various crossover methods, such as single-point crossover, multiple-point crossover, and uniform crossover, may be used.
- Mutation: in order to avoid becoming trapped in a local optimal solution, a mutation process is used.
   In this process, some genes in the chromosomes are selected randomly and their values changed.
   Mutation is generally damaging rather than beneficial to the optimization process. However, it reduces
   the probability of trapping in a local optimum point.
- **Evaluation:** the purpose of this step is to evaluate the goodness of each chromosome. The evaluation is done by finding the objective function value (in this paper, the value of  $Z_1 Z_2$ . The less Z1 and Z2 the better the chromosome is.
- Stopping criterion: There are several strategies for stopping the evolution process of GA. Usually,
   two procedures are adopted as convergence criterion: (1) the iteration stops when the variation in the
   fitness level among generations is within a user-defined range; and (2) when the number of generations

reaches to a predetermined level. Both approaches are applied in this research; while the procedure stop when each stopping criterion pleased first. In this research the predefined range for variation in the fitness level is set 0.001 and the number of generations is set as 100\*the number of decision variables  $(Z_1 \text{ has } 6 \text{ and } Z_2 \text{ has } 5 \text{ decision variables}).$ 

5 In comparison with TELUM and G-LUM, applying the Genetic Algorithms would increase the 6 flexibility of the model, by eliminating the need of testing different initial points. In addition, adding more 7 variables in the model would increase the chance of trapping in a local optimum solution and makes finding 8 the optimum solution harder. In problems with the high local optimum solutions, Genetic Algorithms 9 reduce the chance of trapping in local optimum solution (*27*).

10 11

1 2

3

4

## DATA REQUIREMENT AND DATA PREPARATION

12 The proposed land use model needs six sets of input data. Households, employment, house 13 conditions (total and vacant houses), amount of land in five land use classes, and travel time. These data 14 sets are needed for two periods of time with a time interval of five years. The household data, along with 15 categories (total population, total households, household income, household size, household worker, household seniors, household students, quarter group) collected from census data. This data set is available 16 17 for every 10 years. The employments data containing 20 categories of NAICS codes are available through 18 Longitudinal Employment and Household Dynamics (LEHD). The house condition (the number of total 19 and vacant houses) is collected through census data. The land use condition in five land use classes 20 collected from parcel data (28); Each parcel has a year of the built attribute allowing extending data for 21 previous years. By using this information, land use conditions generated from the year 2000 to 2020 every 22 five years. Lastly, travel time data is obtained from TSTM.

23

## 24 RESULTS

25 This section discusses the result of model validation and accuracy. In order to illustrate the model 26 applicability and validity, LS-LUM and LS-LUMS-WHC are implemented in the state of Tennessee, United 27 States. The state of Tennessee has 95 counties and 3,293 TAZs. Due to data collection limitation (especially 28 parcel data, even though available but need to be requested), in this paper, the model is applied in 39 29 counties (see Figure 2). The selected study area has 1,451 TAZs with a population of 2,881,195 and total 30 employment of 1,755,491 in 2010. To test the model performance, households are modeled in 9 categories 31 (total population, total households, households with 1 to 6 persons, and households with 7 or more persons) 32 and employments modeled by 21 categories (total employment and 20 NAICS employment categories). 33 The three-step validation process is discussed below.

34 First, models are developed for the base year (2010), then backcasting and forecasting accuracies 35 are presented. At each step, the goodness-of-fit measure  $R^2$  and the error, Mean Absolute Percentage Error 36 (MAPE) are provided. The proposed model results are compared with G-LUM. Generally, a model with 37 higher  $R^2$  and smaller MAPE is a better model. In addition, based on Chin (29) study which proposed a rule of thumb for acceptable  $R^2$ , where  $R^2$  greater than 0.66 is substantial, between 0.33 and 0.66 is moderate, 38 39 and less than 0.33 is week, in this paper,  $R^2$  greater than 0.66 is considered as acceptable. Two points 40 should be mentioned. First, since G-LUM does not model total houses, agricultural, and vacant land, the  $R^2$ 41 and MAPE for these variables are provided only for LS-LUM. Second, because the only difference between 42 LS-LUM and LS-LUM-WHOHC is in modeling households, the results for LS-LUM-WHOHC are 43 provided just for households.



Figure 2 The State of Tennessee with 95 Counties and 3293 TAZs; The Model Implemented on the
 Blue Part (39 Counties with 1451 TAZs)

#### 4 Developing the Model for The Base Year 2010

5 In the first step, the model developed for the year 2010. and 2005 is considered as the lag year. As 6 Figures 3 to 8 present, both LS-LUM and G-LUM are fitted very well in all the categories, except for 7 employment NAICS sector 51. The goodness of fit in the households and land use conditions (Figures 3 8 and 5) is better in LS-LUM in comparison with G-LUM. The differences in the land use section are more 9 significant. In addition, Figure 3 shows that removing the *HC's* variables from the model is reduced the  $R^2$ 10 of households, specifically in households with 7 or more person group. Also, the *MAPE* is increased

11 significantly in comparison with both LS-LUM and G-LUM.





Figure 3 The  $R^2$  and *MAPE* of three models for households for the base year 2010





Figure 4 The  $R^2$  and *MAPE* of employments for the base year 2010





#### Figure 5 The R<sup>2</sup> and MAPE of total houses and land use conditions for the base year 2010

#### 3 Model Backcasting Validation

After developing the model for the base year 2010, the model backcast the households, 4 5 employment, total houses, and land use condition in the year 2005. Figures 6 to 9 show backcasting 6 validation results illustrating LS-LUM could predict all categories with acceptable accuracy. However, the 7 performance of LS-LUM in household and land use condition prediction is better than the employment 8 section. The  $R^2$  in all categories is greater than 0.85. In addition, the difference between LS-LUM and G-9 LUM is more significant in these sections. In the employment section, both models show similar goodness 10 of fit (Figure 7). However, the value of MAPE in LS-LUM is lower in all employments' categories in 11 comparison with G-LUM (Figure 7).





Figure 6 The  $R^2$  of backcasting the households for the year 2005





Figure 7 The  $R^2$  and *MAPE* of backcasting employments for the year 2005



1 2

Figure 8 The  $R^2$  and *MAPE* of backcasting land use conditions for the year 2005

# 3 Model Forecasting Validation

Forecasting validation had limits. Since the latest available household data was for the year 2010, the model forecasting accuracy is provided for employments and the land use condition. **Figure 10** shows the  $R^2$  and *MAPE* of the estimated employment in year 2015. Except NAICS sector 53 and 48-49, LS-LUM forecast the employment better compared to G-LUM. Both models have deficits in forecasting NACIS 54 and 71. However, the improvement of LS-LUM in NAICS 11, 21,51, and 56 is significant.

![](_page_16_Figure_1.jpeg)

## 2 Figure 10 The $R^2$ and *MAPE* of forecasting employments for the year 2015

3 Figure 11 presents the models' accuracy in forecasting the land use condition in year 2015. In addition, Figure 12 shows the accuracy of forecasting for the year 2020. Since the parcel data was available 4 5 for the year 2020, it was possible to calculate the goodness of fit of the land use condition in 2020. 6 Generally, similar to backcasting the accuracy of LS-LUM in predicting land use condition is much better 7 than G-LIM. Since the components of LC directly affect the prediction of households and employments 8 and due to the cumulative nature of errors, the accuracy of prediction in this section becomes more 9 important. Moreover, as Figure 12 shows, by increasing the year of prediction, the difference between the 10 goodness of fit in LS-LUM and G-LUM increases. These results show that LS-LUM would work better in 11 a long-range forecasting year.

![](_page_17_Figure_1.jpeg)

1 2

Figure 11 The *R*<sup>2</sup> and *MAPE* of forecasting land Use condition for the year 2015

![](_page_17_Figure_4.jpeg)

![](_page_17_Figure_5.jpeg)

Figure 12 The *R*<sup>2</sup> and *MAPE* of forecasting land use condition for the year 2020

## 5 **DISCUSSION**

6 The overall results showed that, LS-LUM provides better accuracy compared to the similar land 7 use model (G-LUM). Several approaches were tested to improve accuracy. The first action was involving 8 land use conditions in the Employment section directly. In the previous models, the result of the land 9 consumption section (LANCON in TELUM and LUDENSITY in G-LUM) does not directly affect the 10 employment allocation section. In LS-LUM the amount of commercial, industrial, and agricultural land, 11 were added to the *Employment* section. Although significant improvement was not observed, minor 1 improvement in some categories is not deniable. Comparing the results of LS-LUM and LS-LUM-WOHC 2 showed that when House Condition section (*HC*) is added to the model, the model accuracy and stability 3 increased significantly. LS-LUM showed higher  $R^2$  lower and *MAPE* in comparison to LS-LUM-WOHC 4 both in developing and backcasting sections.

5 The new formulation for modeling land consumption shows improved accuracy. LC compared to 6 LANCON in G-LUM, provides higher  $R^2$  and lower MAPE in addition to predicting agricultural and vacant 7 land. The improvements in land consumption section are important from a different point of view; due to 8 the LC results affecting the other sections directly, the model can retain its accuracy for additional years of 9 forecasting.

# 10 CONCLUSIONS

11 The purpose of this paper was to develop a land use model that can be applied to large-scale 12 geographies with reasonable computational time and acceptable accuracy using only open source data. The 13 proposed model, large scale land use model (LS-LUM) incorporates the underlying concepts TELUM and G-LUM with improved model formulation, and enhanced solution algorithm. The improved model 14 15 formulation consists of new variables addition in the form of total and vacant houses. LS-LUM involves 16 the amount of commercial, industrial, and agricultural land in predicting the number of employments. A 17 new evolutionary computation-based solution approach is presented to enhance accuracy and optimality. Although the model shows acceptable results in the form of improved R<sup>2</sup> and lower MAPE for household 18 19 and land type, additional research is needed to enhance the accuracy of the EMP-AL and socioeconomic 20 conditions of large-sale cases.

Future studies can consider the effect of other components, similar to land price and salary. Conducting policy and scenario analysis is another avenue of future research. Finally, improving the calibration accuracy and runtime is another direction for future studies by exploring additional heuristic and other evolutionary algorithms.

25

#### 26 ACKNOWLEDGMENTS

This research was funded by Tennessee Department of Transportation (TDOT). The authors would like to thank TDOT for providing data and support during the research period. In addition, computational facilities at University of Memphis is greatly acknowledged. Any opinions, findings, and conclusions or recommendations expressed are those of the authors and do not necessarily reflect the views of the abovementioned agencies

32

## 33 AUTHOR CONTRIBUTIONS

34 The authors confirm contribution to the paper as follows: study conception and design: AR, SM; data

35 collection: AR, SM, JE; analysis and interpretation of results: AR, SM; draft manuscript preparation: AR,

36 SM. All authors reviewed the results and approved the final version of the manuscript.

#### REFERENCES

- 1. Lowry, I. S. A Model of Metropolis. Rand Corp Santa Monica Calif, 1964.
- 2. Acheampong, R. A., and E. A. Silva. Land Use–Transport Interaction Modeling: A Review of the Literature and Future Research Directions. *Journal of Transport and Land use*, Vol. 8, No. 3, 2015, pp. 11–38.
- Salvini, P., and E. J. Miller. ILUTE: An Operational Prototype of a Comprehensive Microsimulation Model of Urban Systems. *Networks and spatial economics*, Vol. 5, No. 2, 2005, pp. 217–234.
- 4. Moeckel, R., C. Schürmann, and M. Wegener. Microsimulation of Urban Land Use. 2002.
- 5. Waddell, P. UrbanSim: Modeling Urban Development for Land Use, Transportation, and Environmental Planning. *Journal of the American planning association*, Vol. 68, No. 3, 2002, pp. 297–314.
- 6. Hunt, J. D., J. E. Abraham, and D. D. Silva. PECAS—for Spatial Economic Modeling. *Calgary, Alberta: HBA Specto Incorporated*, 2009.
- 7. Martinez, F. MUSSA: Land Use Model for Santiago City. *Transportation Research Record*, Vol. 1552, No. 1, 1996, pp. 126–134.
- 8. Waddell, P. Integrated Land Use and Transportation Planning and Modelling: Addressing Challenges in Research and Practice. *Transport Reviews*, Vol. 31, No. 2, 2011, pp. 209–229.
- 9. Moeckel, R., C. L. Garcia, A. T. M. Chou, and M. B. Okrah. Trends in Integrated Land-Use/Transport Modeling. *Journal of Transport and Land Use*, Vol. 11, No. 1, 2018, pp. 463–476.
- 10. Manual, U. TELUM (Transportation Economic and Land Use Model) Version 5.0.
- 11. Gao, S., E. Lehmer, Y. Wang, M. McCoy, R. A. Johnston, J. Abraham, and J. D. Hunt. Developing California Integrated Land Use/Transportation Model. 2009.
- 12. Abraham, J. E., and J. D. Hunt. Design and Application of the PECAS Land Use Modeling System. *University of California, Davis*, 2003.
- 13. Hunt, J. D. 1 Desing and Application of the PECAS Land Use Modeling System. 2003.
- 14. Spasovic, L. N. TELUM-Interactive Software for Integrated Land Use and Transportation Modeling. New Jersey Institute of Technology. 2013.
- 15. Merlin, L. A., J. Levine, and J. Grengs. Accessibility Analysis for Transportation Projects and Plans. *Transport Policy*, Vol. 69, No. C, 2018, pp. 35–48.
- 16. Dimitrijevic, B. A Method for Assessing Transportation Impacts of New Land Developments Using Integrated Land Use and Transportation Network Modeling.
- 17. Wang, C.-H., N. Chen, and S.-L. Chan. A Gravity Model Integrating High-Speed Rail and Seismic-Hazard Mitigation through Land-Use Planning: Application to California Development. *Habitat international*, Vol. 62, 2017, pp. 51–61.
- 18. Pozoukidou, G. Forecasting Land Use Changes: An Empirical Approach for East Thessaloniki. 2017.
- 19. Morton, B. J. Land Use Forecasting Models for Small Areas in North Carolina. 2013.
- 20. Paul, V. V., and B. Zhou. Documentation for Use of the Gravity-Based Land Use Model (G-LUM): With and Without In-Built Travel Demand Model Assumption. 2009.
- 21. Putman, S. H. Integrated Urban Models: Policy Analysis of Transportation and Land Use. Pion Limited. *London, UK*, 1983.
- 22. Valsaraj, V., K. Kockelman, J. Duthie, and B. Zhou. Forecasting Employment & Population in Texas: An Investigation on TELUM Requirements, Assumptions, and Results, Including a Study of Zone Size Effects, for Austin and Waco Regions.
- 23. Zhou, B., K. M. Kockelman, and J. D. Lemp. Applications of Integrated Transport and Gravity-Based Land Use Models for Policy Analysis. *Transportation research record*, Vol. 2133, No. 1, 2009, pp. 123–132.
- 24. RSG. Tennessee Statewide Travel Model (Version 3) Develpment and Validation Report. 2014.

- 25. Bernardin Jr, V. L., N. Ferdous, H. Sadrsadat, S. Trevino, and C.-C. Chen. Integration of National Long-Distance Passenger Travel Demand Model with Tennessee Statewide Model and Calibration to Big Data. *Transportation Research Record*, Vol. 2653, No. 1, 2017, pp. 75–81.
- 26. Mathew, T. V, S. Khasnabis, and S. Mishra. Optimal Resource Allocation among Transit Agencies for Fleet Management. *Transportation Research Part A: Policy and Practice*, Vol. 44, No. 6, 2010, pp. 418–432.
- 27. Durand, N., and J.-M. Alliot. A Combined Nelder-Mead Simplex and Genetic Algorithm. 1999.
- 28. Tennessee Comptroller of the Threasury. https://comptroller.tn.gov/office-functions/pa/gisredistricting/land-use-data.html.
- 29. Chin, W. W. The Partial Least Squares Approach to Structural Equation Modeling. *Modern methods for business research*, Vol. 295, No. 2, 1998, pp. 295–336.