



**WAGENINGEN**  
UNIVERSITY & RESEARCH

# Mapping and predicting the bio-feedstock availability for material transition

Final wildcard project report

Xuezhen Guo  
Marcel Kornelis  
Daoud Urdu

22-12-2022

This project has been funded by the investment theme Transformative Bioeconomies:  
Towards a materials transition that phases out fossil feedstock

## Introducing the format

When submitting your Wildcard project you committed to providing several deliverables:

1. A short accessible document for the inter- and transdisciplinary group of people involved in the programme that describes your methodological innovation project / proof of principle project and its rationale;
2. A presentation at a community meeting of the investment theme;
3. A report of the results of your learning journey that describes the key lessons learned about your methodological innovation or proof of principle.
4. Additional deliverables formulated by you as part of the submission, labelled 'Project specific deliverables' in this format.

All Wildcard projects already provided presentations as stipulated under 2. This format then is meant to document deliverables 1, 3 and 4.

In section 2 of the format we ask some additional questions related to possible follow-up.

### 1. A short accessible document (max. 600 words)

In order to produce materials from biomass sources to facilitate the bioeconomy transition, it is important to understand the future development of biomass availability. This study aims to use historical data to develop a predictive time-series model that can project the future biobased-feedstock availability from the major selected crops.

The popularity of time-series models can be (at least partly) explained by the fact that they can already be usefully applied if only the performance variables of interest (biomass production in our case) is available. This is, because time-series models can use the own past of the performance variable as a basis for predicting the future values. Another aspect that may explain its popularity, is that these type of models can relatively easy incorporate the opinion of experts and market observers.

The insights derived from this research will help the policy makers (such as the Common Agricultural Policy, the Sustainable Development Goals and the Paris Climate Agreement) to evaluate the potential of using biobased feedstock to replace the traditional fossil feedstock and to make policy concerning materials transition.

#### Innovative idea and objective

The success of model predictions in practice may depend on what is called *forecast efficiency*. This implies that the model can quickly adopt changes in the markets when they occur, and take them adequately into account when a new prediction is made. In other words, each time when new data becomes available a Decision Support System (DSS), including a forecast efficient model, is able to adopt the relevant information from this new data.

Recently, machine learning techniques have been developed to check for forecast efficiency of predictive models. By applying machine-learning techniques to the building of time-series models, we may develop predictive models with two strong features. Firstly, it may help to capture information in the available data, which was not found if only time-series

techniques were used. Secondly, the model can still capture expert opinion in its functional form, as it remains a time-series model in its basis.

Up to the best of our knowledge, we are the first to consider the use of machine-learning techniques to test for the forecast efficiency of time-series models, within the context of biomass production.

### Relevance to the materials transition in textiles and/or building materials?

To facilitate the materials transition (i.e., replacing the fossil-based materials by the biobased materials), forecasting the availability of the biomasses especially from the major sources (crops) are essential. We believe that *forecast efficiency* is a relevant topic within this context, because the production patterns of these markets have to change in order to facilitate the bioeconomy transition. Consequently, useful predictive models should not only give accurate forecasts of the biomass commodities of interest, but also fast adopt information about potential changing patterns, so that predictions can quickly be revised.

### What did you do?

- We downloaded the FAOSTAT data on production and land use for cotton, flax and wood pulps for the five continents (Africa, America, Asia, Europe, Oceania), as well as for the world.
- We developed AutoRegressive (AR) models to predict each biomass variable of interest (a total of 36 models)
- We developed Vector AutoRegressive (VAR) models to predict each biomass variable of interest (a total of 6 models)
- We organized a workshop with WUR researchers where we discussed our approach
- We compared the results of the single equation AR models with those of the VAR model with respect to model accuracy and efficiency.
- Develop the prototype of the dashboard to visualize the results.

### Key deliverables:

- Overview of data sources dataset
- A predictive model, in EViews and R, which applies the AR, VAR, and random forest techniques to be potentially developed into a DSS tool
- A draft manuscript to be submitted to a scientific journal
- Wireframe dashboard as a prototype to visualise results

### Main result, achievement and highlight

Describe the key results of your work. What insights have been generated? What is it you want to highlight?

A draft manuscript to be submitted to a scientific journal, is for the methodology and empirical application developed within this project. This chapter contains preliminary results that are part of research paper to be submitted. It ends with a descriptive section that presents workshop results as part of this study.

A Vector AutoRegressive (VAR) system forms the central part of our approach. The methodology is related to the studies of. VAR models can easily capture (i) both short- and long-run components, (ii) existing dynamics in the relationships between the variables, and (iii) do not require firm prior knowledge about the nature of the causal relationships. VAR models can be specified in levels, in differences and in Error-Correction form, depending on the outcomes of preliminary unit-root and cointegration tests. The most general specification is the Vector Error Correction Model (VECM). For a two-variable case, this model can be written as:

$$\begin{bmatrix} \Delta x_t \\ \Delta y_t \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} t + \sum_{s=1}^{S-1} \begin{bmatrix} b_{s,1} \\ b_{s,2} \end{bmatrix} d_{st} + \sum_{j=1}^J \begin{bmatrix} \pi_{11}^j & \pi_{12}^j \\ \pi_{21}^j & \pi_{22}^j \end{bmatrix} \begin{bmatrix} \Delta x_{t-j} \\ \Delta y_{t-j} \end{bmatrix} + \begin{bmatrix} \alpha_1 & 0 \\ 0 & \alpha_2 \end{bmatrix} \begin{bmatrix} e_{1,t-1} \\ e_{2,t-1} \end{bmatrix} + \begin{bmatrix} u_{1,t} \\ u_{2,t} \end{bmatrix}, \quad (6.1)$$

where  $t = 1, \dots, T$ ,  $x_t$  is the advertising expenditures variable,  $y_t$  is the macro-economic indicator, and  $\Delta$  the difference operator, e.g.  $\Delta x_t \equiv x_t - x_{t-1}$  (see also equation (3.4)).

Deterministic components include the intercept terms ( $\mu_1$  and  $\mu_2$ ), a deterministic trend term ( $t$ ) and seasonal dummy variables ( $d_{st}$ ). The *short-run* inter-relationships between  $x_t$  and  $y_t$  are captured in two ways: through the autoregressive parameters  $\pi_{kl}^j$  ( $k, l = 1, 2$ ,  $j = 1, \dots, J$ ) and through the error-correction terms  $\alpha_k e_{k,t-1}$  ( $k = 1, 2$ ). The former capture the traditional short-run dynamics, with the diagonal elements measuring the impact of own past behavior on current behavior and the off-diagonal elements capturing the lagged cross-effects. The error-correction terms capture the short-run adjustment for temporary deviations from a long-run co-movement or equilibrium between  $x_t$  and  $y_t$ . Such a situation exists when  $x_t$  and  $y_t$  are cointegrated, and reflects the *long-run* inter-relationships between  $x_t$  and  $y_t$ . Finally,  $u_{1,t}$  and  $u_{2,t}$  are multivariate normal disturbance terms with variance-covariance matrix  $\Sigma$ . Instantaneous relationships between  $x_t$  and  $y_t$  are not captured directly in the model, but are reflected in the off-diagonal elements of  $\Sigma$ . When the variables have a unit-root but are not cointegrated, the error-correction parameters  $\alpha_k$  ( $k = 1, 2$ ) become zero. When no unit-root is found, the difference operator is omitted, and the corresponding variable is specified in the levels. In that case, the variable is said to be either level or trend stationary. In the latter case, the deterministic trend component is needed in the specification, while the parameter  $\beta_k$  ( $k = 1, 2$ ) becomes zero in the former case.

### 3.1 Random forest to evaluate the joint forecast efficiency of different models

The joint forecast efficiency is an indicator to reflect whether a forecasting model has captured all the information from the information set (i.e., the explanatory variables of the model). For an efficient model, there should be no residual information existing in information set, which means there is no structural relationship between explanatory variables and the error terms of the forecasting model on the testing dataset.

To quantify the joint forecast efficiency, traditional linear regression model can be replaced by the random forest model (Behrens, et al., 2018) because it can better capture the non-linear relationship between the error terms of forecasting and the explanatory variables. The detailed steps for applying the random forest process are as follows:

- 1) Create and train 100 random forests which contains 1, 2, 3, 4...100 decision trees using the information set of the original bootstrapped data. Using the 100 trained random forest models to predict the error terms of the 100 out-of-bag dataset and derive 100 distances measuring the difference between the “real error terms” and the “predicted error terms”. We then use the median of the 100 derived distances as the distance indicator for the current round. We will repeat 1000 rounds to get a distribution of the median distances.
- 2) We follow the same steps to estimate the random forests 1000 times on a permuted matrix of forecast errors, where the permuted data are computed by sampling without replacement from the original data. We then get another distribution of the median distances.
- 3) We then test if the mean of the median-distance distribution derived based on the original data is significantly smaller than that of the median-distance distribution derived based on the original data. The hypotheses are as follows:

$H_0: \Delta M \geq 0$ ,  $H_1: \Delta M < 0$

$\Delta M$  is the difference of the means of the two distributions.

For more details, please refer to the work of Behrens, et al. (2018).

#### *Empirical application*

##### *4.1 The production of cotton, flex, and roundwood*

We consider the production of cotton, flux, and wood pulp at different levels of aggregation. We employ annual data on production and used area for the period from 1961 to 2020, as collected by FAOSTAT.

Our methodology is applied on the data from 1961 to 2014, and the subsequent six months are used as hold-out sample (i.e. the forecast period accounts for 10% of the sample length).

In addition to the production series, we also analyse area use (measured in hectares), as they may provide a more detailed picture of the underlying dynamics at work. Indeed, when used in combination with the production findings, they may give some important insights on what happened with the production per hectare.

We take logarithms of the variables to reduce potential heteroscedasticity, as is common in recent time-series studies. In general, the first difference of the log-transformed series is a good measure for the growth rate of the original variable (Franses & Koop, 1998).

Unit root tests

We first established the stationarity of the series of interest through unit-root tests. The outcomes of the Augmented Dickey Fuller test for the series are given in Tables 1A to 1G. The procedure of Perron (1989) is used to determine the maximum number of lags ( $J$ ). We start with an equation which includes a deterministic time trend. If the test indicates a stochastic trend, or if the deterministic trend function in a stationary model is found to be insignificant (at a 10% level), we investigate whether the series of interest is level stationary. In our empirical analysis, the cotton-seed case differs in two ways from the flax and wood-pulp cases. First, besides the production and land use for cotton seed, we also analyse the global demand for cotton, chemical fabrics, and wool. For the flax and the wood-pulp markets, we do not have demand observations available. Second, in December 1991, the Soviet Union was dissolved. As a result, the number of cotton-producing countries in Asia increased, whereas the European cotton-seed production decreased. This structural change was not a consequence of changing demand, but the result of a rearrangement in geographical scope. Not dealing with this change may bias the outcome of the unit-root tests. Therefore we included structural break dummies on the level and deterministic trend function in the unit-root tests to allow for this change. As date of change we opted for the year of 1992. The outcomes of the unit-root tests are given in Tables 1A to 1G in Chapter 4 as appendix *Statistical test outcomes*.

#### Forecast efficiency

Concerning the joint forecast efficiency, t-statistics can be used to rank model forecast efficiency. The t-statistics are even addable and their summation of them can even evaluate the overall forecast efficiency of multiple models. See the reference from (Behrens 2020):

“The rank of a forecaster is determined by summing up the t-statistics from the model specifications. The intuition behind this is, that a small negative or positive t-statistic is evidence against a rejection of joint forecast efficiency, whereas a large negative t-statistic leads to a rejection of joint forecast efficiency. Hence, the sum of a given forecaster’s t-statistics over all three scenarios is an indicator of the strength of evidence against joint forecast efficiency of the said forecaster. Summing up the t-statistics implies that a positive t-statistic in one specification can compensate the effect of a negative t-statistic in another specification.”

The results in the tables below show that none of the forecasting models are efficient because all the p values of the t test are significant, which means the random forest derived median distances based on the original dataset are significantly smaller than the median distances based on the permuted dataset. This means there are residual information in the information set that has not been captured by the model.

The outcomes of the forecast efficiency are given in the tables in Chapter 4 as appendix *Statistical test outcomes*.

## Workshop results

To assess and complement the predictive model, a reflexive workshop was organized. This section presents the results of the workshop that was part of this study. It is important to bear in mind the possible bias in the responses of the attendants. More information on the format of the workshop is found as an additional project specific deliverable in chapter 4.

The most striking result to emerge from the workshop was that parameters for the model are more varying than currently is assumed. It is questionable why forecast accuracy is often considered as the main quality attribute, as the model cannot capture changes that will probably arrive in 5 years. An adequate projection model should be flexible enough to capture these changes. The forecast efficiency should provide an answer for this dilemma. However, there are some things that will be the same, that is part of the model. It has commonly been assumed that broadening the scope of efficiency and accuracy, it might improve the overall output of the model with considering the degree of uncertainty. This could imply changes in the model for selection and development of actual DSS applications. For example, boundaries of 90 % could be presented in a graph to show the this quality attribute. The research would have been more relevant if a wider range of quality attributes had been explored.

Further, it was found that the needs and demand of biomass feedstock might be as important as the supply side to have a more integral approach. This finding is complemented with the consideration of existing biobased and non-biobased (fossil-based) products. The growth of crude oil, a well-known raw material for non-biobased products, seems to have large impact and it might be worthwhile to understand the substitution effect that comes along with this material. The substitution flow could provide insights in loosing or gaining positions for specific products in the same market. This is a rather unexpected result and the associated data should be interpreted with caution because, for example, cotton is not per se substitutional from synthetic since there is always fuels needed. Note, that the established models in our study, due to this expert response, incorporates a dataset describing the demand of cotton per kg/per person in the model.

A reasonable approach to tackle the issue of incompleteness could be to conduct a more explorative study on a specific products market to understand the competitors landscape of biobased and non-biobased products. Since phasing out fossil fuel, it is hard to clarify when a certain product, such as cotton, will be phased out. It would still be beneficial to know what the future growth and demand would be of cotton, and what implications could be derived from this. A possible limitation for this is the dynamics that play in the industrial value chain and on a regional level.

The results derived from the predictive model indicate that the growth paths of bio-based and non-biobased products are not interconnected. Although ideally it should be one market, the model approach it currently as separated markets. This result is in line with those obtained during the workshop. We may not find a substitution effect in the data, since awareness of bio-based products just started. Some key suggestions that are derived from the workshop are listed below, while general statements can be found in the section Key message.



- If energy prices increase then changes happen such as looking for alternatives. On the other hand if prices go down, then there are less incentives to produce bio-based products. This is also related to the interconnectedness in a broader sense. For example, a war in Ukraine, has effect on both synthetic and bio-based materials.
- Additional study for the building materials could complement since it is a more event based specific market and less structural. However, this is not the same level as textile, which is directly related to bio-based economy.
- It is, and stays, important to consider the general predications about the material transitions in terms of direction of the transition, speed of transition, impact on sustainability and impact on global economy.

### Key message

What is the key message that people working on the materials transition should remember from your project?

- The forecast accuracy differs substantially across biomass commodities and geographical regions. It is of interest to investigate what the underlying factors of these differences are. Various outcomes may be of interest. For example, it may be the case that region- and commodity specific drivers play an important role in the over-time development of the biomass commodities. It may also be the case, that actors in different regions respond differently to the same global drivers. In both cases, the use of forecast accuracy as a model-building criterion leads to more insights in specific underlying short- and long-run dynamics of biomass commodities.
- The outcomes of the forecast efficiency tests show that both the single equation (AR) and multiple equations (VAR) models do not incorporate all the relevant information from the data that is available. This implies that it is possible to increase insights in the short- and long-run dynamics of biomass commodities with the datasets which are already available. The challenge here is, however, that the tests indicate that this is the case, but it does not give guidelines where to find the relevant information in the data. It may, for example, be related to the interrelationships between various biomass production trends, or to the response behavior in the associated markets. In both cases, the use of forecast efficiency as a model-building criterion insights in specific underlying short- and long-run dynamics of biomass commodities.
- The relative simple AR models are not generally outperformed by the relative complicated VAR models with respect to forecast accuracy and efficiency. This implies that relative parsimonious models can be used to gain insights in underlying short- and long-run dynamics of biomass commodities.
- Focus for this study is on materials, however there is interconnectedness of the energy market and production of materials, such as nylon which is a substitute of flax and cotton. It is important to have non-carbon based energy. However, a link to the energy domain is missing. For example, wood could also be burned to produce energy. There might be existing data from FAO to reuse on half products, such as panel wood.
- In the longer run, transition from a policy perspective is part of the solution, depending on incidental events. It is suggested to prepare some scenario's on big

events and make conclusions of those that are linked to cotton production and sustainability. For example, EU rules that are being developed on climate could bring the market more close. On a national level, in the Netherlands for example, nitrogen crisis and could be an interesting case to analyses due to this finding. A Dutch farmers association, BO Akkerbouw, have brought forward the idea for additional crops that are not so burdensome for the soil to use for biomass. The market for the material transition is however on the international level.

- A substitution matrix could provide insights in the substitutes and minimize, for example, demand for cotton. This approach is followed by an on-going project. A more disruptive measure could be to invest more in the cotton field and tackle barriers. For by-products to use in another domain, it is important to notice that there are qualitative criteria on properties of materials such as a building may not collapse .

Visual abstract

Please place a visual abstract in this box. The box can be bigger or smaller. Please add a caption of the visual abstract below the image.

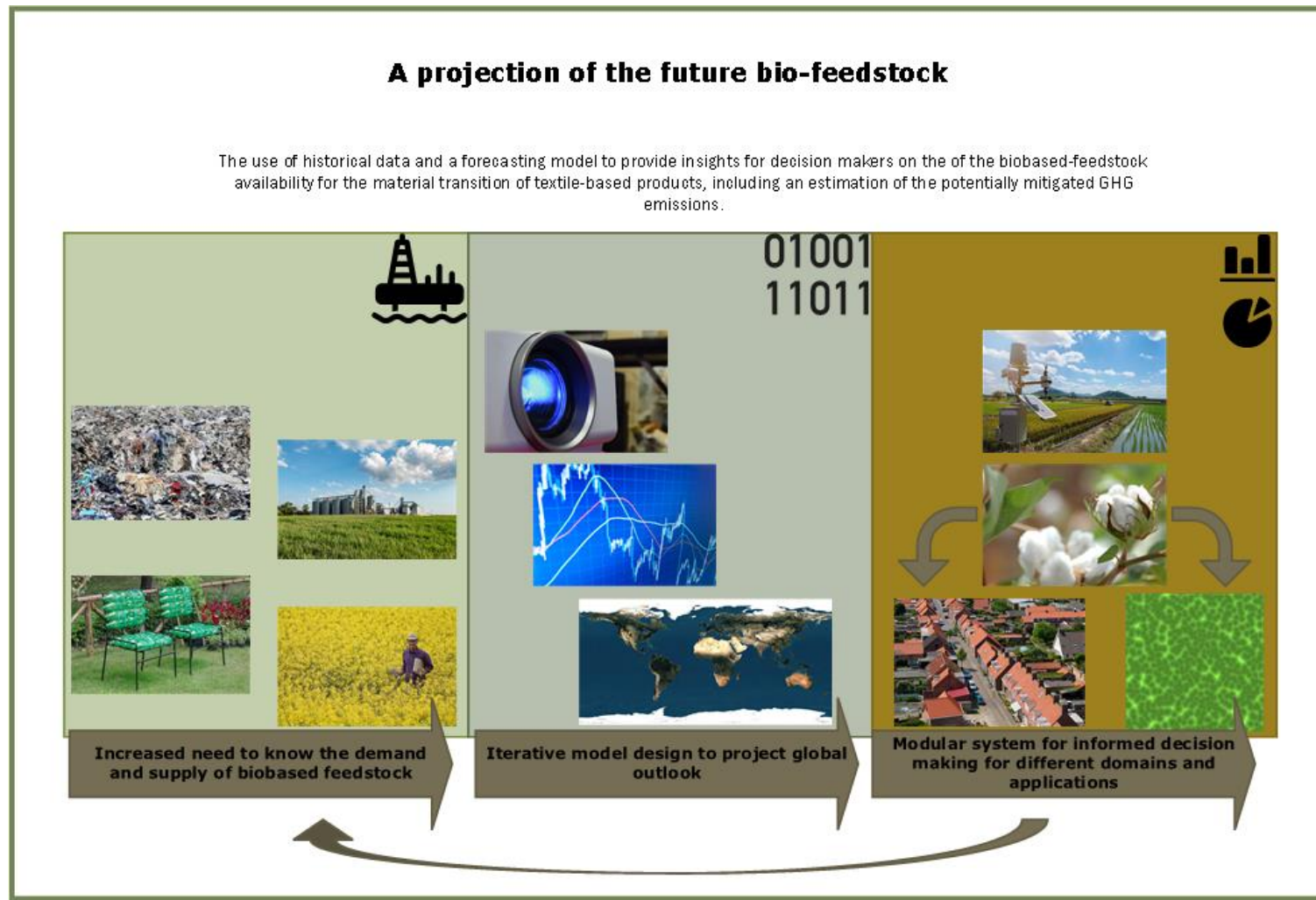


Figure 1: Visual abstract of a projection of future bio-feedstock

## 2. Questions about 'readiness' and possible follow-up (max 200 word)

This section serves the investment theme to understand the development the project has undergone. We aim at selecting Wildcard projects to be taken up by one of the domain flagships (building materials, textiles). To make a selection, we need to know what the progress has been, where the project is now, and what potential there is.

### Where you started

Explain where the project started. E.g. was there already some foundation, or did you have to start from zero?

There was no foundation when the project started, so we started from zero.

### Where are you now

Compared to where you took off, where are you now? What progress has been made? What remains to be done when looking at where you intended to be with this project at the start?

A robust model with certain level of quality attributes. With a focus on efficiency

Remains to be done:

- Finishing the scientific paper
- Algorithm in Python
- Data harmonization: what additional open data sources could be semantically integrated?

### Potential and next steps

How do you currently assess the potential of your project to contribute to the materials transition? What are logical next steps to take it further?

Similar to the key messages given in this report, the following could contribute to potential next steps:

- The forecast accuracy differs substantially across biomass commodities and geographical regions. It is of interest to investigate what the underlying factors of these differences are.
- The outcomes of the forecast efficiency tests show that both the single equation (AR) and multiple equations (VAR) models do not incorporate all the relevant information from the data that is available.
- Generally, in order to answer complex questions within the material transition, an integrated approach is needed. This project could contribute to comprehensive decision making for different actors in the bioeconomy value chain, such as researchers, policy makers, investment decision makers, etc. by facilitating an interoperable infrastructures for applying novel data science techniques. Within this project the process of decision making was narrowed down to time-series models

including specific resources and materials, like cotton, while the infrastructure should be deployable to other problem solving processes.

### Innovation readiness

Where does the project/innovation stand in terms of *readiness*? Is this something that can be piloted or rolled out in the outside world, or is this something that needs some further development and (lab) testing before it can be piloted in society? Is it possible/meaningful to indicate an ‘innovation readiness’ level using the below scale? If so, how would you score your project idea?

Innovation readiness score	Innovation readiness level	Description
0	Idea	Genesis of the innovation. Formulating an idea that an innovation can meet specific goal.
1	Hypothesis	Conceptual validation of the idea that an innovation can meet specific goals and development of a hypothesis about the initial idea.
2	Basic Model (unproven)	Researching the hypothesis that the innovation can meet specific goals using existing basic science evidence.
3	Basic Model (proven)	Validation of principles that the innovation can meet specific goals using existing basic science evidence.
4	Application Model (unproven)	Researching the capacity of the innovation to meet specific goals using existing applied-science-evidence.
5	Application Model (proven)	Validation of the capacity of the innovation to meet specific goals using existing applied science evidence.
6	Application (unproven)	Testing of the capacity of the innovation to meet specific goals within a controlled environment that reflects the specific spatial-temporal context in which the innovation is to contribute to achieving impact.
7	Application (proven)	Validation of the capacity of the innovation to meet specific goals within a controlled environment that reflects the specific spatial-temporal context in which the innovation is to contribute to achieving impact.
8	Incubation	Testing the capacity of the innovation to meet specific goals or impact in natural/real/uncontrolled conditions in the specific spatial-temporal context in which the innovation is to contribute to achieving impact with support from an R&D.
9	Ready	Validation of the capacity of the innovation to meet specific goals or impact in natural/real/uncontrolled conditions in the specific spatial-temporal context in which the innovation is to contribute to achieving impact without support from an R&D.

Table 1: Innovation readiness levels as distinguished by Sartas et al, 2020.

We think that this project stands in 3.

### 3. Learning Journey (max 300 words)

We would like to understand a bit more about the process you went through, and whether and how being part of the investment theme Transformative Bioeconomies influenced your learning. We ask the project leaders to consult others when answering these questions.

1. Did your Wildcard project involve new collaboration with disciplines or people? If so, briefly explain what was new.

New collaboration in terms of teams forming, meaning working for the first time with known colleagues. Also in terms of interdisciplinarity, exploring and discussing new types methods for new context and different backgrounds.

2. If applicable, did the new collaboration alter your original thinking about the topic? Did it change research directions or courses of action? If so, briefly characterize how.

It did alter some original thinking, such as the use of existing models like Magnet.

3. Did interactions during community days and/or meetings organized by the investment theme alter your original thinking about the topic? Did such interactions change research directions or courses of action? If so, briefly characterize how.

The collaborations within the investment theme opened new domains of research and ideas to contribute to sustainability in general.

4. Did you meet any challenges during implementation of your wildcard project? If so, what kind of challenges were these?

Workload from other projects requested challenges in project management. Furthermore, in the beginning lack of specific expertise since a former project member left. This is, fortunately, solved quickly with the replacement of a colleague.

5. If applicable, how were these challenges eventually addressed? Did activities organized by the investment theme contribute to overcoming challenges? If so, briefly indicate how.

From both of the mentioned challenges, it did not contribute specifically.

6. Has your involvement in the investment theme resulted in any new initiatives or spin-offs that would probably not have emerged if you had not participated? If so, briefly indicate how these new initiatives came about.

There are no new initiatives resulted as part of our involvement.

#### 4. Additional project specific deliverables

Copy-paste the deliverables provided in your submission document and explain how you have met these deliverables. If deliverables could not be reached, please explain.

##### Additional deliverables proposed when submitting the Wildcard project

Copy/paste from proposal

To enable producing materials from biomass main and side streams, the pre-requisite is to understand the availability of the biomass at different levels (global, regional, national). This insight will help the policy/decision makers to evaluate the potential of using biobased feedstock to replace the traditional fossil feedstock. Moreover, to make long-term policy/decisions, the policy/decision makers do not only need to know the current situation of the biomass availability but also the long-term trend of the development in the future. This is especially relevant when long term investment and intervention strategies needed to be made by the organization such as CGIAR, IFAD, etc. For this sake, predictive modelling using the historical data to project the future biobased feedstock potentials is desired. Traditionally, the predictive models in this field usually just applied simple (linear) extrapolation, which is not adequate to explore the hidden/complex correlations between different variables. To solve this problem, machine learning models can be applied. Finally,

it is very interesting for the stakeholders to know what are the mitigated GHG emissions because of the switch of feedstock sources from fossil to biobased ones. This proposal analyses the supply part of the potential available bio-based material which can contribute to both of the domain flagships, textiles and building materials.

Based on the aforementioned reasons, in this project, we would like to:

- 1) Mapping the current availability of the biomass main and side streams at global, regional, national levels using public databases (e.g., FAOSTAT) and literature
- 2) Predicting the availability of the biomass feedstock in the future (e.g., 2035) using the historical data coupled with other relevant data (e.g., population, land use changes) with machine learning. To make the results more robust, we will consult domain experts to develop multiple scenarios that represent the maximum, average and minimum scenarios.
- 3) Quantifying the mitigated GHG emissions because of the transition (this will be conducted on the high level due to the limited budget)

Scientific relevance:

To the best knowledge of us, this is the first study which aims to quantify/forecast the availability of the biomass and the potential mitigated GHG emission using AI (or machine learning) at global, regional and national levels. It therefore fill the knowledge gap in the stream of literature concerning biomass availability calculation.

This idea develops the novel option of mapping and predicting the potential of the biomass feedstocks for producing biobased materials as well as the potential impact on GHG emission reduction. It therefore has a big potential of contributing to material transition.

Societal value and relevance

This results of the project can add value to policy and societal decision making.

Policy/decision makers can use the insights of the biomass feedstock availabilities as well as its implications on GHG emission reduction to accordingly make investment/intervention decisions to accelerate material transition. It also contributes to the development of theory of change in the field of bio-materials transition.

Activities & deliverables:

Task 1: Connect different data sources to create the database. Calculate the current availability of the biomass main and side streams at global, regional, national levels

Task 2: Develop the machine learning model to predict the availability of the biomass feedstock in the future. Take into account existing models, such as MAGNET.

Task 3: A workshop to elicit expert-options to calibrate the model with different scenarios (max, average, min)

Task 4: Calculate the mitigated GHG emissions due to the transitions

Task 5: Develop the prototype of the dashboard to visualize the results

Status of each project specific deliverable  
Please report the status of each deliverable.

*Statistical test outcomes*

**Table XX Statistical test outcomes for cotton seed production (in tonnes)**

Variable	t-value	lag	Outcome
Global level	-5.50 <sup>a</sup>	1	trend stationary
Continental level			
Africa	-3.73 <sup>a</sup>	0	trend stationary
Americas	-5.73 <sup>a</sup>	1	trend stationary
Asia	-3.63 <sup>a</sup>	3	trend stationary
Europe	-.84	0	unstable
Oceania	-2.01	0	unstable

<sup>a</sup> significant at the 5% level, for which the critical value is  $-3.50$  (Enders 1995 p. 419).

We first established the stationarity of the series of interest through Augmented Dickey-Fuller unit-root tests. The outcomes are given in Table XX, and indicate that cotton production at the global level, and in the series of Africa, Americas, and Asia are trend stationary ( $p < 0.05$ ), while the unit-root null hypothesis was not rejected for the Europa and Asia series.

**Table XX Statistical test outcomes for cotton seed land use (in ha)**

Variable	t-value	lag	Outcome
Global level	-5.79 <sup>b</sup>	1	level stationary
Continental level			
Africa	-2.87	1	unstable
Americas	-4.23 <sup>a</sup>	0	trend stationary
Asia	-2.93	3	unstable
Europe	-2.18	0	unstable
Oceania	-1.48	0	unstable



<sup>a</sup> significant at the 5% level, for which the critical value is  $-3.50$  (Enders 1995 p. 419). <sup>b</sup> significant at the 5% level, for which the critical value is  $-2.93$  (Enders 1995 p. 419).

Asia series.

**Table XX Statistical test outcomes for flax production (in tonnes)**

Variable	t-value	lag	Outcome
Global level	-6.23 <sup>a</sup>	4	trend stationary
Continental level			
Africa	-2.03	0	unstable
Americas	.032	0	unstable
Asia	-2.31	0	unstable
Europe	-3.30 <sup>b</sup>	0	level stationary
Oceania	n.a.	n.a.	n.a.

<sup>a</sup> significant at the 5% level, for which the critical value is  $-3.50$  (Enders 1995 p. 419). <sup>b</sup> significant at the 5% level, for which the critical value is  $-2.93$  (Enders 1995 p. 419).

**Table XX Statistical test outcomes for roundwood production (in tonnes)**

Variable	t-value	lag	Outcome
Global level	-1.74	4	unstable
Continental level			
Africa	-1.73	4	unstable
Americas	-2.31	4	unstable
Asia	-2.32	3	unstable
Europe	-2.17	1	unstable
Oceania	-.30	0	unstable

<sup>a</sup> significant at the 5% level, for which the critical value is  $-3.50$  (Enders 1995 p. 419). <sup>b</sup> significant at the 5% level, for which the critical value is  $-2.93$  (Enders 1995 p. 419).

#### 4.3. Cointegration tests

Stationary or non-stationary series enter our VAR framework in levels or first differences, respectively. We use Akaike's Information Criterion to derive the optimal number of lags in the model. If there are two or more non-stationary series in the model, we first established a cointegrating relationship between the non-stationary series of interest by means of a

Johansen trace test. In the cointegrating relationship, we follow (reference needed) in that we allow for an intercept but not for a trend in the cointegrating equation, and that there is no deterministic trend in the data. In the case of cotton production in Europe and Oceania, no cointegrating relationship was found between these two series. We therefore, consider a mixed VAR model for the continental level, and not an VECM.

**Table XX Cointegration test outcomes**

Variable	Hyp. number of CE(s)	trace-test value	lag timing break(s)	Outcome
Continental level				
Europe Oceania	and none	14.18	1	No cointegration
	at most one	3.69	1	No cointegration

**Table XX Cointegration test outcomes**

Variable	Hyp. number of CE(s)	trace-test value	lag timing break(s)	Outcome
Continental level				
Africa, Europe Oceania	Asia, none <sup>a</sup> and	54.43	1	One cointegrating relationship
	at most one	25.69	1	
	at most two	12.16	1	
	at most three	4.12	1	

<sup>a</sup> significant at the 5% level, for which the critical value is 54.08 (MacKinnon et al. 1999).

#### 4.4. Forecast performance

We now concentrate on the forecast performance of the alternative models. The results are presented in Table XX. Based upon the RMSE, we conclude that the VAR model shows a better forecast performance in comparison to the single equation models.

**Table XX Forecast performance test outcomes Cotton production in tonnes**

Variable	RMSE	Efficiency
Single equation models		
Global level	.102	
Continental level		
Africa	.095	
Americas	.201	
Asia	.250	
Europe	.150	
Oceania	1.102	
VAR/VECM models		
Africa	.093	
Americas	.176	
Asia	.227	
Europe	.272	
Oceania	.639	

**Table XX Forecast performance test outcomes Cotton land use in ha**

Variable	RMSE	Efficiency
Single equation models		
Global level	.07	
Continental level		
Africa	.05	
Americas	.13	
Asia	.11	
Europe	.15	
Oceania	.95	
VAR/VECM models		
Africa	.06	

Americas	.19
Asia	.09
Europe	.14
Oceania	.88

---

**Table XX Forecast performance test outcomes flax production in tonnes**

Variable	RMSE	Efficiency
Single equation models		
Global level	1.89	
Continental level		
Africa	.006	
Americas	.038	
Asia	1.165	
Europe	.389	
Oceania	n.a.	
VAR/VECM models		
Africa		
Americas		
Asia		
Europe		
Oceania	n.a.	

---

**Table XX Forecast performance test outcomes roundwood production in tonnes**

Variable	RMSE	Efficiency
Single equation models		
Global level	.023	
Continental level		
Africa	.019	

Americas	.037
Asia	.017
Europe	.065
Oceania	.066
VAR/VECM models	
Africa	
Americas	
Asia	
Europe	
Oceania	n.a.

---

The outcomes of the unit-root tests are given in Tables 1A to 1G.

**Table 1A. Statistical test outcomes for cotton seed production (in tonnes)**

Variable	t-value	lag	Outcome
Global level	-5.22 <sup>a</sup>	1	trend stationary
Continental level			
Africa	-2.20	0	unstable
Americas	-6.80 <sup>a</sup>	1	trend stationary
Asia	-.21 <sup>c</sup>	2	unstable
Europe	-14.15 <sup>c</sup>	4	trend stationary with break
Oceania	-1.81	0	unstable

<sup>a</sup> significant at the 5% level, for which the critical value is -3.50 (Enders 1995 p. 419).

<sup>b</sup> significant at the 5% level, for which the critical value is -2.93 (Enders 1995 p. 419).

<sup>c</sup> we applied known-breakpoint unit root test, for which the critical value at the 5% level is -4.24 (Perron 1994, p.135).

**Table 1B Statistical test outcomes for the global demand for cotton, wool, and chemical processing (in tonnes)**

Variable	t-value	lag	Outcome
Cotton	-4.12	3	trend stationary
Wool	-0.26	2	unstable

Chemical -0.61 0 unstable

<sup>a</sup> significant at the 5% level, for which the critical value is  $-3.50$  (Enders 1995 p. 419).

<sup>b</sup> significant at the 5% level, for which the critical value is  $-2.93$  (Enders 1995 p. 419).

<sup>c</sup> we applied known-breakpoint unit root test, for which the critical value at the 5% level is  $-4.24$  (Perron 1994, p.135).

**Table 1C Statistical test outcomes for cotton seed land use (in ha)**

Variable	t-value	lag	Outcome
Global level	-5.57 <sup>b</sup>	1	level stationary
Africa	-2.75	1	unstable
Americas	-.60	3	unstable
Asia	-4.96 <sup>c</sup>	3	trend stationary with break
Europe	-25.04 <sup>c</sup>	4	trend stationary with break
Oceania	-1.67	3	unstable

<sup>a</sup> significant at the 5% level, for which the critical value is -3.50 (Enders 1995 p. 419).

<sup>b</sup> significant at the 5% level, for which the critical value is -2.93 (Enders 1995 p. 419).

<sup>c</sup> we applied known-breakpoint unit root test, for which the critical value at the 5% level is -4.24 (Perron 1994, p.135).

**Table 1D Statistical test outcomes for flax land use (in ha)**

Variable	t-value	lag	Outcome
Global level	.27	3	unstable
Africa	-2.95 <sup>a</sup>	0	trend stationary
Americas	-2.83	3	unstable
Asia	-17.43 <sup>c</sup>	0	trend stationary with break
Europe	-157.13 <sup>c</sup>	0	trend stationary with break
Oceania	n.a.	n.a.	n.a.

<sup>a</sup> significant at the 5% level, for which the critical value is -3.50 (Enders 1995 p. 419).

<sup>b</sup> significant at the 5% level, for which the critical value is -2.93 (Enders 1995 p. 419).

<sup>c</sup> we applied known-breakpoint unit root test, for which the critical value at the 5% level is -4.24 (Perron 1994, p.135).

**Table 1E Statistical test outcomes for wood pulp land use (in ha)**

Variable	t-value	lag	Outcome
Global level	-2.75	0	unstable
Africa	-2.90 <sup>a</sup>	3	trend stationary
Americas	-2.83	3	unstable
Asia	-.80	0	unstable
Europe	-.70	0	unstable
Oceania	.46	0	unstable

<sup>a</sup> significant at the 5% level, for which the critical value is -3.50 (Enders 1995 p. 419).

<sup>b</sup> significant at the 5% level, for which the critical value is -2.93 (Enders 1995 p. 419).

<sup>c</sup> we applied known-breakpoint unit root test, for which the critical value at the 5% level is -4.24 (Perron 1994, p.135).

**Table 1F Statistical test outcomes for flax production (in tonnes)**

Variable	t-value	lag	Outcome
Global level	-6.23 <sup>a</sup>	4	trend stationary
Continental level			
Africa	-2.01	0	unstable
Americas	-1.44	0	unstable
Asia	-4.12 <sup>a</sup>	4	trend stationary
Europe	-3.25 <sup>b</sup>	0	level stationary
Oceania	n.a.	n.a.	n.a.

<sup>a</sup> significant at the 5% level, for which the critical value is  $-3.50$  (Enders 1995 p. 419).

<sup>b</sup> significant at the 5% level, for which the critical value is  $-2.93$  (Enders 1995 p. 419).

<sup>c</sup> we applied known-breakpoint unit root test, for which the critical value at the 5% level is  $-4.24$  (Perron 1994, p.135).

**Table 1G Statistical test outcomes for wood pulp production (in tonnes)**

Variable	t-value	lag	Outcome
Global level	-3.79 <sup>b</sup>	2	level stationary
Continental level			
Africa	-2.10	0	unstable
Americas	-3.65 <sup>b</sup>	2	level stationary
Asia	-1.53	2	unstable
Europe	-2.95 <sup>b</sup>	0	level stationary
Oceania	-2.23	3	unstable

<sup>a</sup> significant at the 5% level, for which the critical value is  $-3.50$  (Enders 1995 p. 419).

<sup>b</sup> significant at the 5% level, for which the critical value is  $-2.93$  (Enders 1995 p. 419).

<sup>c</sup> we applied known-breakpoint unit root test, for which the critical value at the 5% level is  $-4.24$  (Perron 1994, p.135).

Stationary or non-stationary series enter our AR and VAR frameworks in levels or first differences, respectively. We use Akaike's Information Criterion to derive the optimal number of lags in the model.

#### 4.4. Forecast performance

We now concentrate on the forecast performance of the alternative models. The results are presented in Table 2A to 2G. Based upon the RMSE, we conclude that the VAR model shows a better forecast performance in comparison to the single equation models.

**Table 2A Forecast accuracy outcomes for cotton-seed production in tonnes**

Variable	Root Mean Squared Error	
	AR model	VAR model
Global level	.09	
Continental level		
Africa	.09	.16
Americas	.17	.26
Asia	.11	.17
Europe	.20	.20



Oceania	.67	.94
---------	-----	-----

**Table 2B Forecast accuracy outcomes for the global demand for cotton, wool, and chemical processing (in tonnes)**

Variable	Root Mean Squared Error	
	AR model	VAR model
Cotton	.05	.03
Wool	.04	.06
Chemical	.06	.04

**Table 2C Forecast accuracy outcomes for cotton seed land use (in ha)**

Variable	Root Mean Squared Error	
	AR model	VAR model
Global level	.19	
Continental level		
Africa	.16	.36
Americas	.37	.27
Asia	.28	.26
Europe	.38	.31
Oceania	1.74	1.61

**Table 2D Forecast accuracy outcomes for flax land use (in ha)**

Variable	Root Mean Squared Error	
	AR model	VAR model
Global level	.20	
Continental level		
Africa	.04	.09
Americas	.03	.06
Asia	1.26	1.28
Europe	.39	.42
Oceania	--	--

**Table 2E Forecast accuracy outcomes for wood pulp land use (in ha)**

Variable	Root Mean Squared Error	
	AR model	VAR model
Global level	.01	
Continental level		
Africa	.01	.01
Americas	.01	.01
Asia	.02	.02
Europe	2.89	.05
Oceania	3.39	.11

**Table 2F Forecast accuracy outcomes for flax production (in tonnes)**

Variable	Root Mean Squared Error
----------	-------------------------

	AR model	VAR model
Global level	.12	
Continental level		
Africa	.01	.06
Americas	.02	.02
Asia	.23	.45
Europe	.18	.21
Oceania	--	--

**Table 2G Forecast accuracy outcomes for wood pulp production (in tonnes)**

Variable	Root Mean Squared Error	
	AR model	VAR model
Global level	.07	
Continental level		
Africa	.13	.08
Americas	.06	.05
Asia	.03	.04
Europe	.05	.04
Oceania	.14	.04

Table XX Forecast efficiency outcomes for cotton yield

Cotton_yie Id	Single-equation model			Var model		
	Beta	t value	P value	Beta	t value	P value
				-		
Africa	-0.0116815	27.52	<2e-16 ***	0.015401 9	24.33	<2e-16 ***
				-		
Americas	-0.0225551	23.84	<2e-16 ***	0.001247 4	8.478	2.25e-13 ***
				-		
Asia	-0.0012248	-10.4	<2e-16 ***	-0.053984	24.67	<2e-16 ***
				-		
Europe	-0.0033498	16.33	<2e-16 ***	-0.030512	19.73	<2e-16 ***
				-		
Oceania	-1.1697	44.13	<2e-16 ***	-0.96019	21.62	<2e-16 ***

		-	<2e-16
World	-0.0014026	9.993	***

Total t scores for the single equation model: -122.22

Total t scores for the VAR model: -98.828

Table XX Forecast efficiency outcomes for cotton land use

Cotton_land_ use	Single-equation model	t		Beta	t	
		value	P value		value	P value
	Beta	-	<2e-16		-	<2e-16
Africa	-0.0208725	86.46	***	-0.07861	43.01	***
		-	<2e-16		-	<2e-16
Americas	-0.024613	21.01	***	0.002844	9	13.66
		-	<2e-16		-	<2e-16
Asia	-0.0209072	24.79	***	0.006685	8	17.42
		-	<2e-16		-	<2e-16
Europe	-0.0102793	42.62	***	0.003829	5	24.98
		-	<2e-16		-	<2e-16
Oceania	-1.67559	72.37	***	-0.61546	27.35	***
		-	<2e-16		-	<2e-16
World	-0.0200829	45.63	***			

Total t scores for the single equation model: -247.25

Total t scores for the VAR model: -126.42

Table XX Forecast efficiency outcomes for flax yield

Flax_yiel d	Single-equation model	Var model
----------------	--------------------------	--------------

	<b>Beta</b>	<b>t value</b>	<b>P value</b>	<b>Beta</b>	<b>t value</b>	<b>P value</b>
Africa	-7.84E-06	-10.52	<2e-16 ***	-1.74E-03	-19.56	<2e-16 ***
Americas	-5.47E-05	-12.11	<2e-16 ***	-3.48E-04	-37.44	<2e-16 ***
Asia	-0.83049	-33.39	<2e-16 ***	0.010712	-5.469	3.4e-07 ***
Europe	-0.016411	-32.59	<2e-16 ***	-0.04388	-30.76	<2e-16 ***
World	-0.049875	-42.4	<2e-16 ***			

Total t scores for the single equation model: -88.61

Total t scores for the VAR model: -93.229

Table XX Forecast efficiency outcomes for flax land use

<b>Flax_land_use</b>	<b>Single-equation model</b>			<b>Var model</b>		
	<b>Beta</b>	<b>t value</b>	<b>P value</b>	<b>Beta</b>	<b>t value</b>	<b>P value</b>
Africa	-1.18E-05	-8.96	2.02e-14 ***	-2.39E-03	35.17	<2e-16 ***
Americas	-2.89E-05	10.15	-	-8.48E-05	16.24	<2e-16 ***
Asia	-0.015925	12.39	<2e-16 ***	-0.29807	22.06	<2e-16 ***
Europe	-3.16E-05	8.651	9.51e-14 ***	-9.79E-05	11.08	<2e-16 ***
World	-4.12E-04	24.79	<2e-16 ***			

Total t scores for the single equation model: -40.151

Total t scores for the VAR model: -84.55

Table XX Forecast efficiency outcomes for wood pulp yield

Woodpulp_yield	Single-equation model			Var model		
	Beta	t value	P value	Beta	t value	P value
Africa	-1.44E-03	-	<2e-16 ***	-3.51E-03	-	<2e-16 ***
Americas	-4.03E-05	-	1.21e-12 ***	-5.23E-04	-28.3	<2e-16 ***
Asia	-8.88E-05	-	<2e-16 ***	-6.59E-04	-	<2e-16 ***
Europe	-4.22E-04	-	<2e-16 ***	-1.77E-04	-	<2e-16 ***
Oceania	-2.16E-03	-	<2e-16 ***	-6.37E-04	-	<2e-16 ***
World	-4.42E-04	-	<2e-16 ***			

Total t scores for the single equation model: -112.249

Total t scores for the VAR model: -183.76

Table XX Forecast efficiency outcomes for wood pulp land use

Woodpulp_land_use	Single-equation model			Var model		
	Beta	t value	P value	Beta	t value	P value
Africa	-6.79E-06	-42.6	<2e-16 ***	-3.12E-06	-	<2e-16 ***
Americas	-9.93E-07	-	<2e-16 ***	-1.16E-05	-26.7	<2e-16 ***
Asia	-8.81E-07	-	<2e-16 ***	-3.67E-06	-	<2e-16 ***

Europe	-7.25E-07	9.131	8.58e-15 ***	-1.92E-05	17.64	<2e-16 ***
Oceania	-9.60E-05	9.608	7.82e-16 ***	-7.78E-04	13.71	<2e-16 ***
World	-1.61E-06	11.56	<2e-16 ***			

Total t scores for the single equation model: -103.419

Total t scores for the VAR model: -86.41

It is not possible to conclude whether or not the VAR model has a better efficiency than the single-equation model (or the other way around) based on the total t scores for forecasting models as presented above.

## Links to or copies of deliverables

Please provide links to or copies of deliverables below. You may insert them as Annexes in this document.

### *Workshop outline*

#### **Title**

The prediction of bio-feedstock availability for material transition as well as the impact on GHG emission mitigation.

#### **Introduction**

To enable producing materials from biomass main and side streams, the pre-requisite is to understand the availability of the biomass at different levels (global, regional, national) in the long term. This insight will help the policy- and decision makers to evaluate the potential of using biobased feedstock to replace the traditional fossil feedstock. This is especially relevant when investment and intervention strategies needed to be made. To this end, predictive models, that use historical data to forecast the future biobased feedstock potentials are desired. In our project we aim to give predictions on availabilities of the biomass materials of interest, such as cotton, flax and wood. So, this work analyses the supply part of the potential available bio-based material which can contribute to both of the domain flagships, textiles and building materials.

#### **Workshop agenda**

- 09:00 – 09.15 Presentation wildcard project – Development of forecast models and goal of Today's workshop by Daoud Urdu and Xuezhen Guo
- 09.15 – 09.50 Preliminary results of different forecast models & explanation of qualitative characteristics by Marcel Kornelis
  - What do we mean by convergence level
  - What do we mean by stable vs. unstable trend
  - What do we mean by unexpected shocks
  - What do we mean by forecasting efficiency
- 09.50 – 10.00 Coffee break
- 10.00 – 10.15 Round 1 – First response of all participants
- 10.15 – 10.30 Round 2 – Vignettes with future scenario's with respect to the qualitative characteristics.
- 10.30 – 10.45 Round 3 – Forecasting efficiency of the two models, which includes judgement on prediction model variables, and prioritization and preferences. In groups of 2.

- 10.45 – 11.00 Round 4 – (Optional or closure) Political issues related to the forecast and reflection as a researcher.

**Date and Time:**

Thursday October 13<sup>th</sup>, 2022. 9 am until 11 am

**Experts invited:**

- Paulien Harmsen - Domain Textile
- Marieke Meeusen – Domain Building materials
- Wolter Elbersen – Materials expert
- Jan Broeze – Materials expert
- Anne-Charlotte Hoes – Transition knowledge
- Bert Annevelink – Biobased economy modelling
- Rene van Ree – EU Call coordinator Bioeconomy
- Harriette Bos – As client representative: main beneficiary
- Geerten Hengeveld – Modeller
- Rens Vliegthart - Modeller

**Wild card team / Facilitators:**

- Marcel Kornelis – Time series analysis modeller
- Xuezhen Guo – Machine Learning modeller
- Daoud Urdu – Information- and project manager

Hybrid form: location and meeting link in the invitation

**Vignettes – supplementary part of the workshop**

Possible future scenarios (or ‘vignettes’) represent ways to understand the roles of different stakeholders with the presented forecast. For example, it could reveal alternative ways the forecast could impact the processes of policy makers. This opens the reflection of Today’s workshop of participants to alternative possibilities. After the presentation, participants are asked to share their intuitive preferences for one or more of these vignettes, and tell their reasons for that choice to the group. In this phase participants are asked to reflect, based on their intuition.



### **Vignette 1 – Accuracy & Efficiency**

I see my competitors shifting to round wood as a product to sell in the building materials domain. As a manager of a company in this domain, I need to have a forecast of what the production volume of this material is on the long-term. This forecast will support me in my decision making whether to invest in this material. More importantly, I need to know if qualitative attributes of this forecast affect trends in prices or production.

### **Vignette 2 – Energy crisis**

As a policy maker for the energy domain there are many considerations to take into account since we are in the middle of an energy crisis. I know that my policy colleagues for the infrastructure & construction and food & agriculture domain are looking for alternative crops that should support the transition to a circular bioeconomy. What is the carbon footprint of these crops? More importantly, how could this transition give sustainable pathways for our energy crisis?

### **Vignette 3 – Geo-political challenges**

Different political challenges could accelerate the trade of specific crops that are more or less related to biomass production. Especially, the relationship of United States and China for the long run is worthwhile to consider when extrapolating results derived from the forecast. Are there any specific events foreseen in this context? Are there any other mentionable relationships of relevant actors?

### **Vignette 4 – Subsidies contributions**

Next month there will be a decisive meeting with my colleagues and partners for subsidies on cotton production for several developing countries. As an advisor, I have the privilege to prepare this meeting and provide supportive evidence on the amount and effectiveness of subsidies. Besides a forecast of cotton production, I would be interested on the impact of cotton prices that will be caused by the granted subsidy from developed countries. What is the role of these main actors in this decision: the coordinating institute, cotton companies and producers?

### **Vignette 5 – Greenhouse Gas Emissions**

The life cycle assessment (LCA) of different product categories are seemingly important for decisions in the agri-food domain that will have a footprint. The calculation of carbon footprints for common goods, such as animals and crops, are well developed. However, which LCA methods could support forecasts of alternative crops, such as flux, cotton and round wood? What are the GHG emissions of these products?

### **Vignette 6 – European sustainability policies**

EU policies on sustainability vary from Farm2Fork to Green Deal and addresses different sustainability facets, such as environmental, climate and social. The policies that mention

circularity and bio economies could be of interest for the forecast of this wild card project. Which underlying values, that are part of these policies, are of importance for the transformative bio economies investment theme?