

Проект нацелен на разработку алгоритма, который бы позволил автоматически находить формулы в древнеисландских сагах. Формула - группа слов, которая употребляется регулярно в одних условиях для выражения некоторой идеи, она может иметь постоянный или вариативный лексический состав (ср. в русских сказках зачин "жили-были"). На основе корпуса из 49 саг был составлен список nграмм, который постепенно фильтровался или кластеризовался с помощью различных эвристик, которые направлены на выделение тех или иных предполагаемых свойств формул (например, похожий контекст с помощью тематического моделирования). С помощью методов анализа естественного языка и анализа данных (векторные семантические модели, разные методы кластеризации, тематическое моделирование) был создан алгоритм, который находит как некоторые часто встречающиеся, так и еще не рассматривавшиеся в научной литературе формулы.