

СОЗДАНИЕ УНИФИЦИРОВАННОЙ НАЦИОНАЛЬНОЙ МЕДИЦИНСКОЙ НОМЕНКЛАТУРЫ

Раузина С.Е., доцент каф. медицинкой кибернетики и информатики им. С.А. Гаспаряна, зав. лаб. семантического анализа медицинской информации ИЦТМ РНИМУ им. Н.И. Пирогова

Международный конгресс ИТМ,
Москва, 12-13 октября 2023





СТРАТЕГИЧЕСКИЙ ПРОЕКТ РНИМУ ИМЕНИ Н.И. ПИРОГОВА «ИНСТИТУТ ЦИФРОВОЙ ТРАНСФОРМАЦИИ МЕДИЦИНЫ»



Программа академического лидерства «Приоритет 2030»

*Стратегический проект РНИМУ имени Н.И. Пирогова
«Институт цифровой трансформации медицины»*



Задачи на разработку:

- *Информационно-поисковая система для предоставления врачу семантически связанной информации*
- *Сервисы поддержки принятия клинических решений*
- *Среда для обучения работе в среде МИС МО студентов «врачебных» факультетов и последипломного образования*
- *Семантический тренажер для подготовки врачей-аналитиков*

Национальная медицинская терминологическая база для унификации работ по созданию поисковых и интеллектуальных систем, а также обработки неструктурированных медицинских текстов



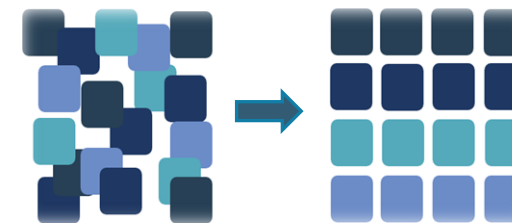
ПРОБЛЕМЫ НЕСТРУКТУРИРОВАННОЙ КЛИНИЧЕСКОЙ ИНФОРМАЦИИ

Отсутствие общепринятой формализации при описании клинической картины пациента, единой медицинской терминологии приводит к:

- Искажению интерпретации информации
- Проблемам анализа на основе данных
- Проблемам создания dataset для проведения научных исследований
- Возможности построения одноплатформенных систем ИИ и информационно-поисковых алгоритмов

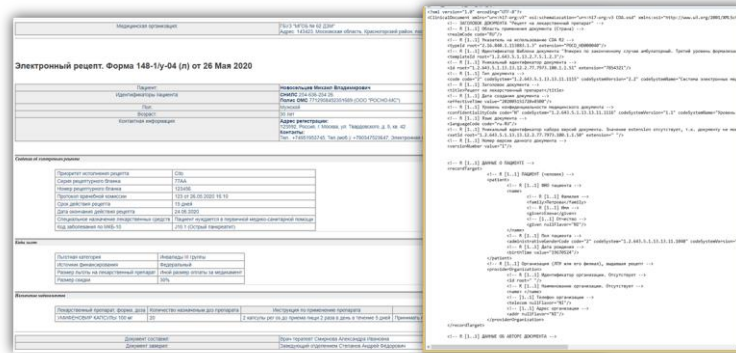


Автоматизация обработки медицинской информации
требует унификации решений при создании алгоритмов ее
поиска и консультативной поддержки врача

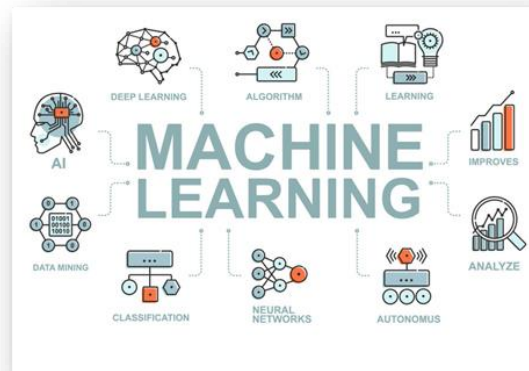


ПУТИ УНИФИКАЦИИ РЕШЕНИЙ

1. Разработка структурированных электронных медицинских документов (СЭМД)



2. Развитие алгоритмов обработки неструктурированных медицинских текстов



Разработка унифицированной медицинской терминологии



Инженеры по знаниям



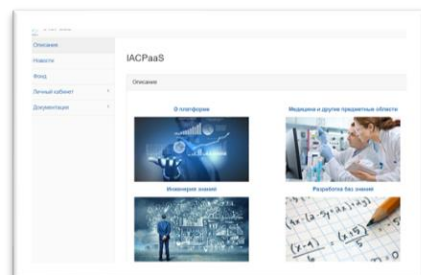
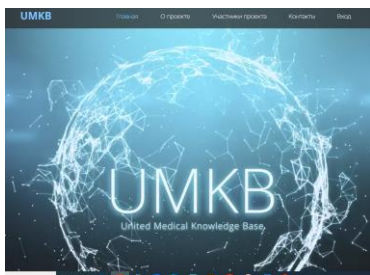
ФУНДАМЕНТ ДЛЯ ПОСТРОЕНИЯ УНИФИЦИРОВАННОЙ НАЦИОНАЛЬНОЙ МЕДИЦИНСКОЙ НОМЕНКЛАТУРЫ (УНМН)

Систематизированный словарь терминов
Связанное представление терминов/концептов
Широкий спектр областей медицины

Отечественный опыт

УМКВ - Объединенная база медицинских знаний

Терминология, созданная на платформе IACPaaS



Система унифицированного медицинского языка, Unified Medical Language System

UMLS

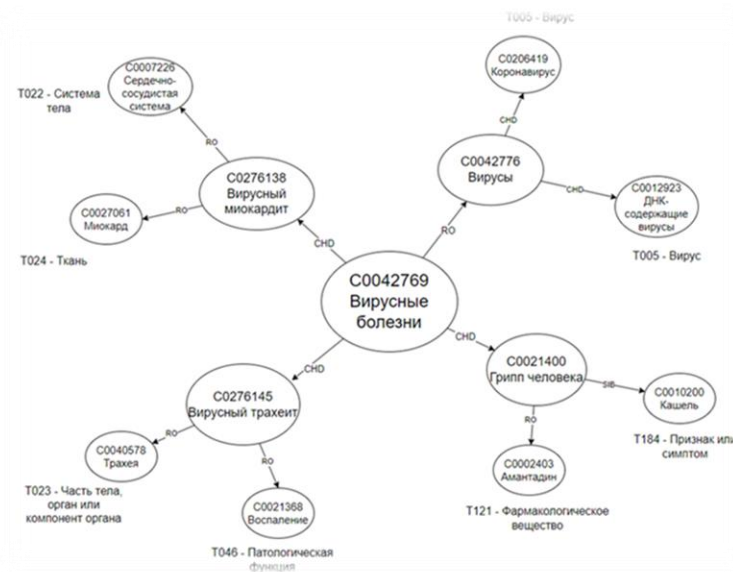
4,6 млн **11,2 млн** **98 млн** **76 актуальн.**

Понятий

Синонимов

Связей

Справочников



SNOMED CT

LOINC
From Regenstrief



MEDICOMP
SYSTEMS

ICD-10-CM

RxNorm

ОЦЕНКА ПОТЕНЦИАЛЬНОЙ ВОЗМОЖНОСТИ ИСПОЛЬЗОВАНИЯ UMLS

Верификация: анализ клинических рекомендаций в части описания клинической картины, данных анамнеза и факторов риска развития заболеваний по **22** группам нозологий.

Всего выделено **579** медицинских терминов

- Аналогичное или полностью схожее наименование – **74%** терминов
- Частичное совпадение – **16%** терминов
- Не найдено - **10%** терминов

Причиной неполного совпадения или ненайденных терминов чаще всего является их сложное с медицинской точки зрения звучание, вероятно, иное в англоязычной среде, требующее привлечения профильных специалистов

№	Наименование КР	ID	Год принят.
1	Аллергический ринит	KP261	2020
2	Острый синусит	KP313	2021
3	Острые респираторные вирусные инфекции у взрослых	KP724	2021
4	Хронический тонзиллит	KP683	2021
5	Острый тонзиллит и фарингит	KP306	2021
6	Острый обструктивный ларингит и эпиглоттит	KP352	2021
7	Паратонзиллярный абсцесс	KP664	2021
8	Грипп у взрослых	Проект	2021
9	Бронхиальная астма	KP359	2021
10	Внебольничная пневмония у взрослых	KP654	2021
11	Хроническая обструктивная болезнь легких	KP603	2021
12	Хронический бронхит	KP655	2021
13	Эмфизема легких	KP656	2021
14	Идиопатический легочный фиброз	KP677	2021
15	Новая коронавирусная инфекция (Covid-19)	Врем. МР	2022
16	Стабильная ишемическая болезнь сердца	KP155	2020
17	Ишемический инсульт и транзиторная ишемическая атака у взрослых	KP171	2021
18	Гастрит и дуоденит	KP708	2021
19	Язвенная болезнь	KP277	2020
20	Хронический панкреатит	KP273	2020
21	Рак желудка	KP574	2020
22	Рак молочной железы	KP379	2021

Например:

*аускультация легких
при эмфиземе*

RUS: коробочный звук

ENG: Hyperresonance

*Геморрагия при
панкреатите*

Симптом Тужилина

СОЗДАНИЕ УНИФИЦИРОВАННОЙ НАЦИОНАЛЬНОЙ МЕДИЦИНСКОЙ НОМЕНКЛАТУРЫ (УНМН)

1 версия

➤ Переводы

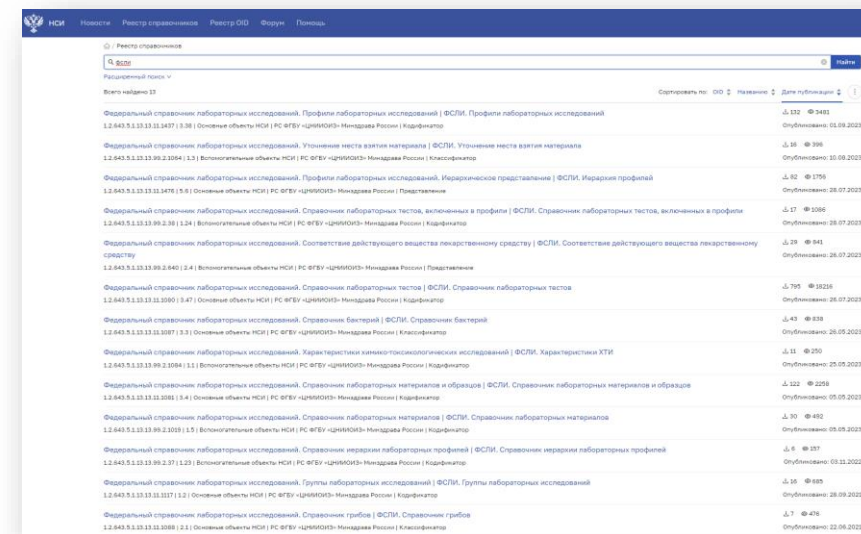
- ✓ Исходно в UMLS представлено **3%** русскоязычных термина
- ✓ Нейросетевые переводы в медицине дают большой % ошибок
- ✓ Переведено **190** тыс. терминов из семантических групп *Симптомы, Клинические находки, Диагностика*



➤ Сопоставление с утвержденными национальными справочниками (ФР НСИ МЗ РФ)

SNOMED_CT, LOINC, RadLex

- ✓ МКБ-10 (100%)
- ✓ ФСЛИ (32%)
- ✓ ФСИДИ (75%)
- ✓ Анатомические локализации (90%)
- ✓ Выявленные патологии (100%)



Всего экспертно верифицированных переводов - **> 215** тыс.
исходно представлено на русском языке в UMLS - **304** тыс.

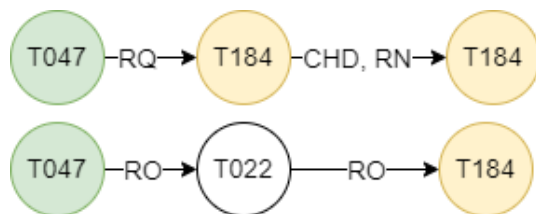
По оценке с помощью разработанной метрики: в UMLS присутствует **≈ 1,5** млн уникальных клинически значимых терминов

МЕТОДЫ АВТОМАТИЗИРОВАННОЙ РАБОТЫ С ТЕРМИНОЛОГИЧЕСКОЙ БАЗОЙ: ИСХОДНЫЕ ЗАКОНОМЕРНОСТИ

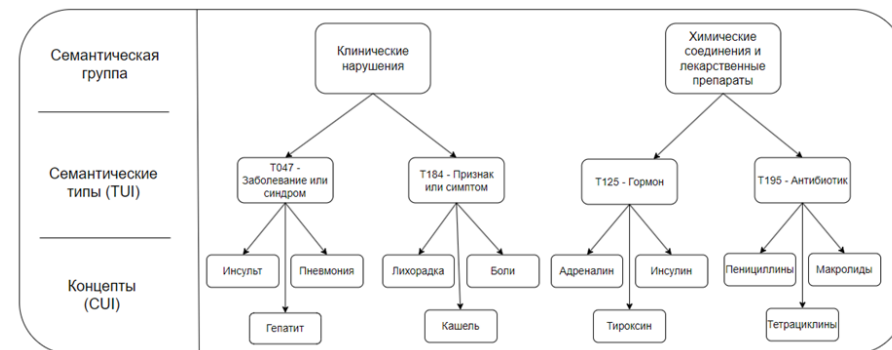
Использование закономерностей в семантической организации терминов и связей, исходно представленных в UMLS

- ✓ Выбор клинически значимых справочников
- ✓ Выявление релевантных групп терминов и типов связей
- ✓ Использование иерархических зависимостей

- ✓ Паттерны для выделения симптомов
- ✓ Паттерны для инструментальной и лабораторной диагностики
- ✗ Паттерны для выбора методов лечения

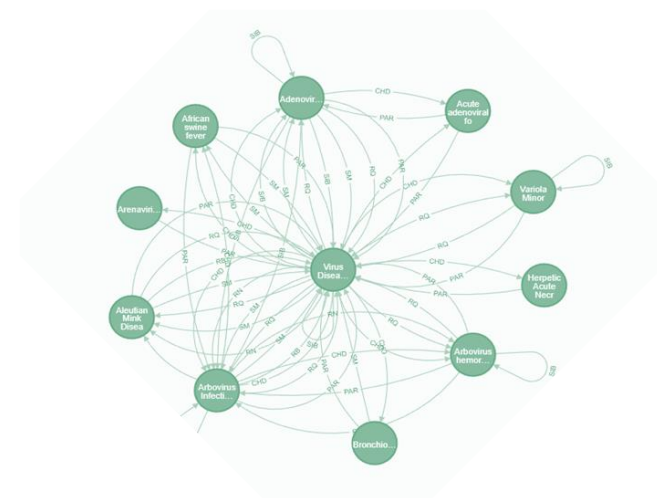


менее 70%
релевантной
информации



127 семантических типов терминов, образующих **15** групп:
(7 социально-экономических и 8 биомедицинских)

Связи объединены в **11** групп и свыше **980** подтипов



МЕТОДЫ АНАЛИЗА СЕМАНТИЧЕСКИХ СЕТЕЙ

✓ Связность концептов:

- прямые связи
- кратчайшие пути
- графовое расстояние
- геометрические фигуры

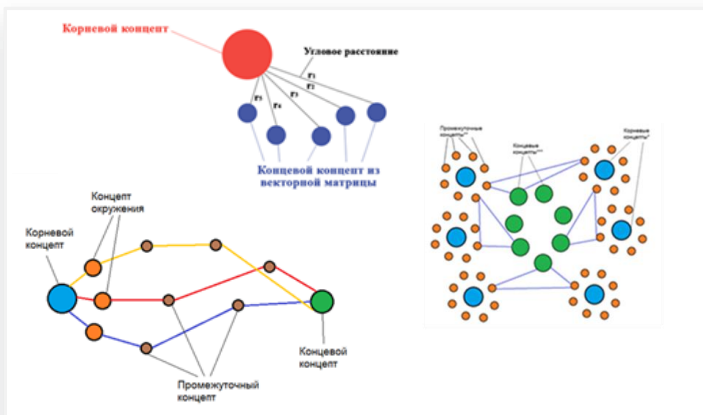


Взвешенный коэффициент кластеризации*:

обеспечивает учет количественной меры связности концептов и степень ее неоднородности

✓ Метрики оценки клинической значимости и специфичности терминов

инструменты ранжирования узлов графа



$$ВКК_n = C_n \cdot R_n^{-1,29}$$

где C_n – число геометрических контуров или незамкнутых путей, в образовании которых участвует узел n ; R_n – число прямых связей между узлом n и любыми другими вершинами графа

Клиническая значимость термина: на основе анализа атрибутов метатезауруса (тематическая группа, актуальность, источник-справочник и др.)

Неспецифичные (обобщающие) термины: «симптом», «признак», «заболевание», «пациент»

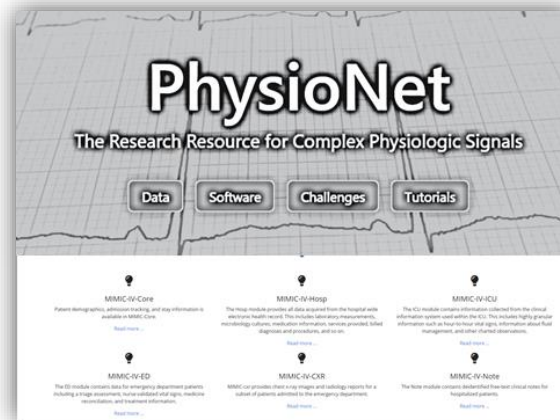
*Астанин П. А., Раузина С. Е., Зарубина Т. В. Построение этиопатогенетического образа концептов метатезауруса UMLS с использованием графовых метрик // Программные системы: теория и приложения. 2023. Т. 14. № 3. С. 59–94. (Рус., англ.). https://psta.psir.ru/2023/3_59-94

ДОБАВЛЕНИЕ ПРЯМЫХ СВЯЗЕЙ В УНМН

✓ 27 млн аннотаций к медицинским англоязычным статьям PubMed



✓ Данные реальной клинической практики: dataset MIMIC-IV (330 тыс. ЭМК)



Лаборатория вычислительной физиологии Массачусетского технологического института
PhysioNet

Семантический анализатор **SemRep** (набор лингвистических правил для анализа англоязычных текстов) \approx 4 млн уникальных связей для УНМН

Агрегатор концептов UMLS **MetaMap** (автоматическое извлечение концептов из англоязычных текстов) \approx 16 млн уникальных связей для УНМН между симптомами и заболеваниями

ПОСТРОЕНИЕ БАЗЫ ЗНАНИЙ НА ОСНОВЕ УНМН

- Закономерности, исходно представленные в UMLS
- Методы анализа семантических сетей
- Метрика автоматизированной оценки клинической значимости концептов
- Метрика оценки относительной специфичности терминов
- Получение дополнительных прямых связей

с **70%** до **90%**
релевантной
информации:

*искомый диагноз
встречается среди
трёх первых
заболеваний в
ранжированном
перечне*

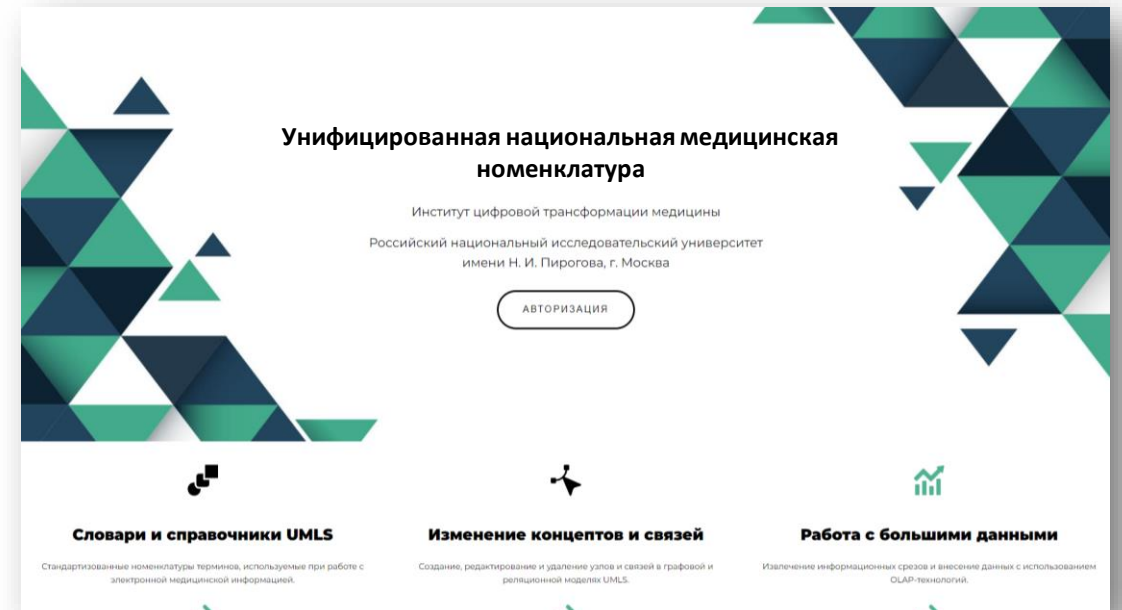
Добавление полученных метрик к УНМН в качестве атрибутов позволило получить онтологическую модель для создания компонентов информационно-поисковых систем и СППКР

ПРОГРАММНЫЕ РЕШЕНИЯ

Аналитическая система для проведения исследовательских работ и обучения

3 основных рабочих блока:

- ✓ Нормативно-справочная информация – международные и национальные словари;
- ✓ Редактор базы знаний – инструмент для наполнения УНМН экспертными знаниями;
- ✓ Интеллектуальное ядро – комплекс аналитических решений, предназначенных для извлечения OLAP-срезов из УНМН.



ПРОГРАММНЫЕ РЕШЕНИЯ

Прототип информационно-поисковой системы (ИПС) для автоматического определения наиболее вероятного класса заболеваний при составлении дифференциально-диагностического ряда произвольного клинического образа и выбора методов диагностики для их уточнения

- ✓ Выделение медицинских понятий из свободно вводимых текстов
- ✓ Поиск отрицаний

Укажите жалобы или симптомы

кашель
постоянный
насморк потеря
обоняния
повышение
температуры

Поиск 🔍

Очистить

C0562483 персистирующий кашель ✖
C0003126 anosmia ✖
C1260880 ринорея ✖
C0015967 гипертермия ✖

Подобрать диагноз

Заболевания **Диагностические методы**

Очистить

C0562483 персистирующий кашель ✖
C0003126 anosmia ✖
C1260880 ринорея ✖
C0015967 гипертермия ✖

Подобрать диагноз

Острый назофарингит
■■■■■ 1 >

Хронический бронхит
■■■■■ 0.892 >

Грипп
■■■■■ 0.88 >

Пневмония вирусная
■■■■■ 0.868 >

- ✓ Вывод ранжированного списка заболеваний на основе введенных симптомов

ОБРАБОТКА НЕСТРУКТУРИРОВАННОГО ТЕКСТА

Шаг 1: Очистка (регулярные выражения), коррекция ошибок (алгоритм Левенштейна)

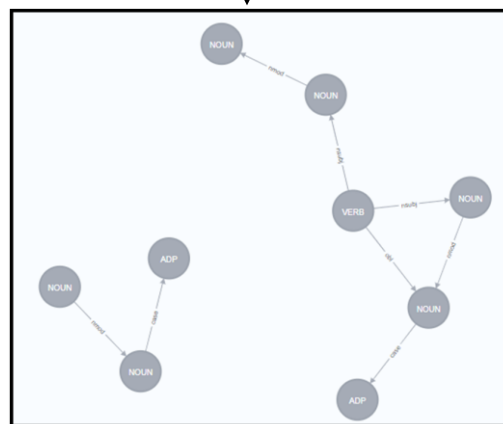
Шаг 2: Лемматизация – приведение к нормальной форме слова (pymorphy2, stanza*)

Шаг 3: Построение дерева синтаксических связей предложения (stanza depparser)*

Боль, оды~~ж~~ка, повышенная температура

боль оды~~ш~~ка повышенная температура

боль одышка повышенный температура



- Боль
- Одышка
- Гипертермия

Поиск построчных совпадений с нормализованной леммой из метатезауруса

Проверка графовых представлений концептов на предмет вхождения в синтаксическое дерево предложения (neo4j)

*нейронная сеть stanza BERT

НАПРАВЛЕНИЯ РАЗВИТИЯ

Разработка новых национальных словарей, например, группы справочников, имеющих отношение к симптомам и клиническим находкам

Использование УНМН для формализации и унификации документов ЭМК, а также при создании СЭМД (единое кодирование информации в ЭМК)

В качестве перспективных проектных задач (в основе которых лежит УНМН) следует назвать необходимость **создания полноценных информационно-поисковых систем**, рассчитанных на пользователя-врача, пациента и исследователя, а также **системы для разработки СППКР**, в которых будут реализованы автоматизированные рабочие места эксперта и инженера по знаниям



БЛАГОДАРЮ ЗА ВНИМАНИЕ

rauzina_se@rsmu.ru