

4. Титов Ю.В., Фарсобина В.В. Сравнительное тестирование авторубрикаторов // «Интеллектуальные информационные технологии. Концепции и инструментарий» Сб. тр. Института системного анализа РАН, 2005. Mode access: <http://www.cognitive.ru/innovation/sbornic7/index.htm>.

5. Косинов Д.И. Использование статистической информации при выявлении схожих документов // Интернет-математика 2007: сб. работ участников конкурса науч. проектов по информ. поиску. Екатеринбург, 2007. С. 205-207.

*Анатолий Юрьевич ОЩЕПКОВ —
аспирант кафедры информационных систем
aoschepkov@gmail.com*

*Александр Григорьевич ИВАШКО —
зав. кафедрой информационных систем,
доктор технических наук, профессор
ivashco@mail.ru*

*Институт математики и компьютерных наук
Тюменский государственный университет*

УДК 519.152

ПОСТРОЕНИЕ АЛГОРИТМА НАХОЖДЕНИЯ «ОПТИМАЛЬНОГО УЗЛА» ДЛЯ РАСПРЕДЕЛЕНИЯ ТРАФИКА «WEB-КОНФЕРЕНЦИИ» В РАСПРЕДЕЛЕННОЙ СИСТЕМЕ

ALGORITHM DEVELOPMENT FOR FINDING THE "OPTIMAL NODE'S" FOR TRAFFIC DISTRIBUTION OF "WEB CONFERENCE" IN CLUSTERING SYSTEM

АННОТАЦИЯ. В статье рассматривается построение алгоритма балансера в кластерной системе для проведения «WEB-конференций», основанного на теории массового обслуживания.

SUMMARY. The article describes the development of algorithm balancer in the cluster system for «Web Conference», based on the theory of waiting lines.

КЛЮЧЕВЫЕ СЛОВА. Веб-конференция; балансер, СМО.

KEY WORDS. Web Conference; Load Balancing; theory of waiting lines.

В последнее время, в связи с внедрением в учебный процесс WEB-2.0 технологий, широкое распространение получили так называемые «вебинары», или «веб-конференции». Среди основных возможностей конференцсвязи можно выделить: проведение слайдовых презентаций; VoIP (аудиосвязь через компьютер в режиме реального времени; видео в режиме реального времени); Whiteboard (электронная доска для комментариев); Screen sharing (удаленный рабочий стол) — совместное использование приложений).

Для университета использование веб-конференций — это возможность для преподавателя работать удаленно так, чтобы студенту казалось, что преподаватель находится рядом. Однако применение этих технологий затруднено в связи с необходимостью использования сервера с высокой производительностью и широкополосного интернет-канала.

Одним из эффективных способов повышения производительности вебинаров является их кластеризация.

В кластерной системе можно выделить следующие группы участников (рис. 1):

- Load Balancing (далее по тексту — балансир) — устройство, служащее для выравнивания нагрузки [2];
- server1...server5 — независимые узлы (с производительностью μ_n) в кластерной системе, между которыми распределяется нагрузка;
- клиент — клиентское приложение, формирующее запросы к системе, общая интенсивность которых составляет λ .

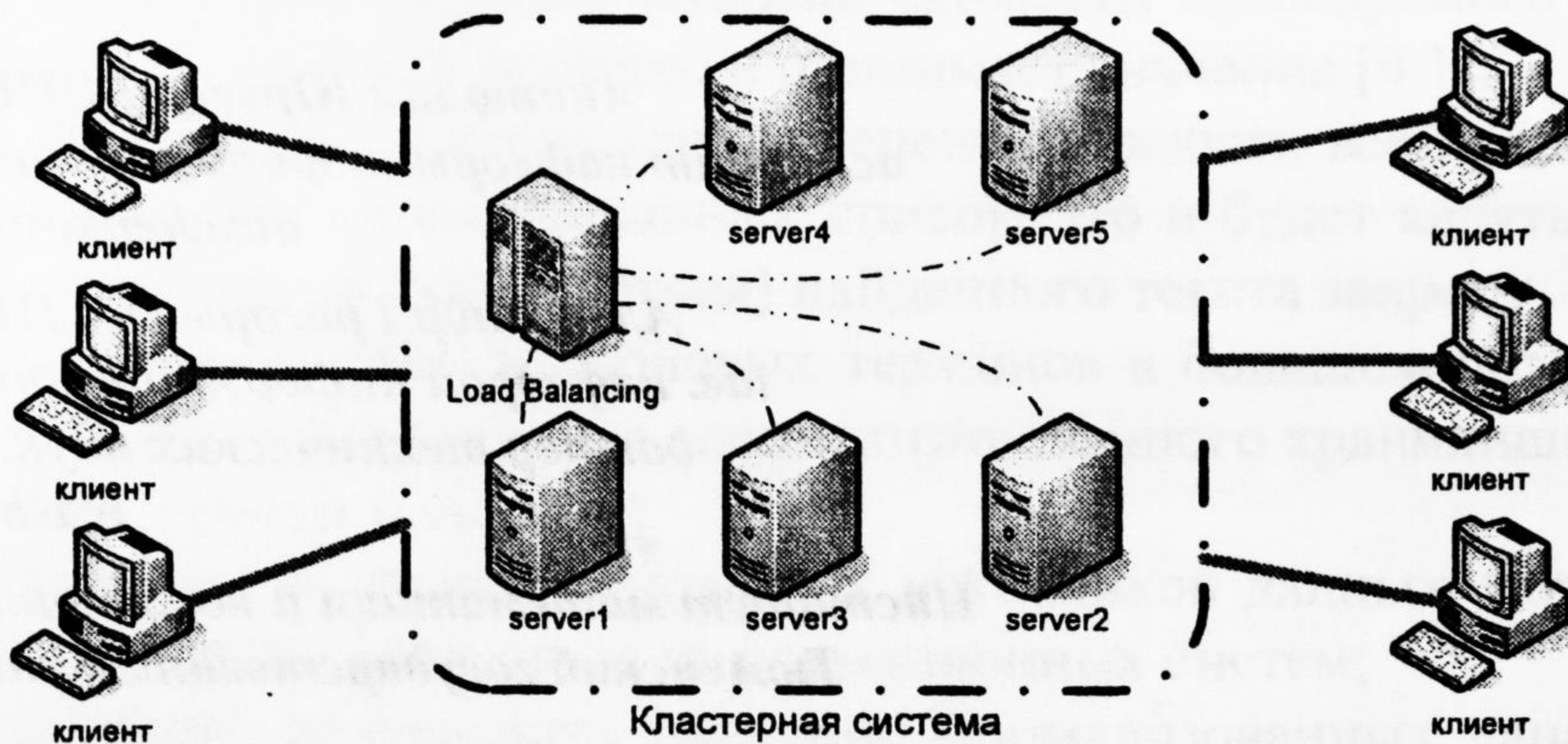


Рис. 1. Кластерная система

Основным элементом в кластерной системе является балансир, основной задачей которого является выбор узла для последующего перенаправления к нему.

Множество состояний балансера можно определить как

$$B = (Y, K), \quad (1)$$

где $K = \{k_j | j \in N\}$ — множество клиентских запросов, поступающих в балансир, которые необходимо перенаправить на узел кластерной системы; Y — множество узлов в кластерной системе.

Клиентский запрос характеризуется двумя параметрами:

$$k_j = (ip_j, id_room), \quad (2)$$

где ip_j — клиентский адрес, с которого сделан запрос; id_room — идентификатор «комнаты», к которой подключается клиент.

Каждый узел в кластерной системе характеризуется следующими параметрами:

$$y_i = (IP_i, Tr_i, P_i, uCount_i), \quad (3)$$

где $uCount$ — динамический параметр — количество пользователей на данном узле; IP — статический параметр — диапазон адресов для клиентов, которые будут считаться локальными по отношению к данному узлу; Tr — динамический параметр — доступный сетевой трафик в Кб/сек; P — динамический параметр — уровень использования памяти.

Работу кластерной системы можно представить в виде системы массового обслуживания с параллельными каналами [3]. При этом она характеризуется количеством, производительностью серверов и их пропускной способностью интернет-каналов; количеством пользователей системы, формирующих запросы. В силу того, что вероятность запроса клиента на подключение к системе в данный момент времени не зависит от запроса в предыдущий момент времени, характер их поступления является пуассоновским процессом, а учи-

тывая то, что вероятность перехода системы из одного состояния в другое определяется только конфигурацией системы в данный момент времени, последовательность состояний системы образуют цепи Маркова [3].

Рассмотрим систему, состоящую из двух серверов *a* и *b* с различной производительностью. При поступлении запроса на вход системы один из каналов немедленно приступает к ее обслуживанию. Время, потраченное на обслуживание запроса, распределено по экспоненциальному закону [5] с параметром μ_a (если заявку обслуживает элемент *a*) и μ_b (если заявку обслуживает элемент *b*).

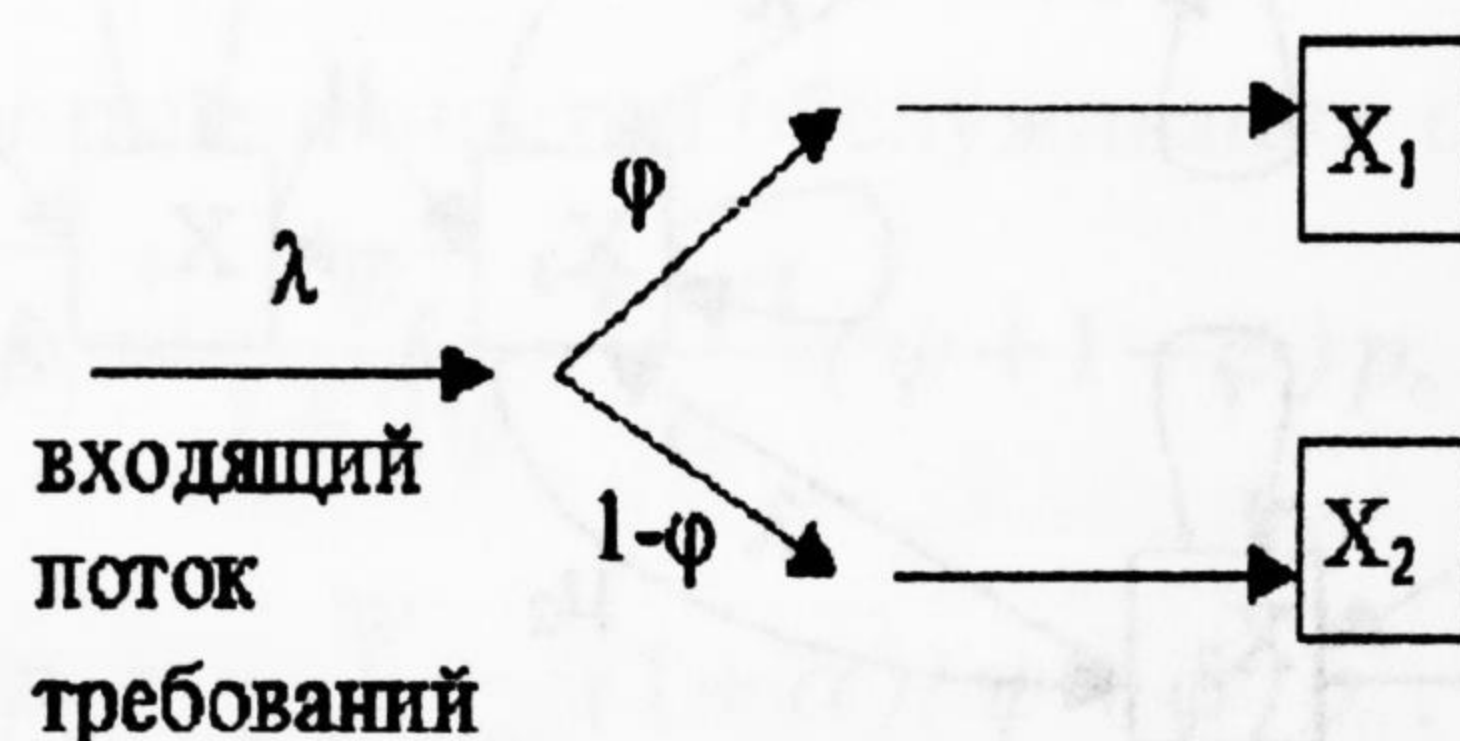


Рис. 3. Двухканальная СМО с параллельными каналами

На рис. 3 показана двухканальная система массового обслуживания с параллельными каналами различной интенсивности обслуживания и вероятностью φ выбора первого канала. Полная интенсивность обслуживания определяется как сумма интенсивностей каждого из узлов. Для распределения потока заявок между каналами обслуживания введем балансирующую вероятность [4]. Требование, которое поступает в момент отсутствия в системе других требований и может поступить на обслуживание в любой прибор, выбирает прибор *i* с вероятностью φ_i . Для двух приборов:

$$\varphi_1 = \varphi \text{ и } \varphi_2 = 1 - \varphi_1, \quad 0 \leq \varphi \leq 1$$

Допустим, что прибор 1 обслуживается быстрее, чем прибор 2, или $\mu_1 > \mu_2$. На основании указанных предположений можно записать следующие условия работы системы:

- при $\varphi = 0$ требование всегда выбирает медленно действующий прибор;
- при $\varphi = 0,5$ требование выбирает случайно один из двух приборов;
- при $\varphi = 1$ требование всегда выбирает быстродействующий прибор.

Возможны и промежуточные значения φ , и тактика обеспечения равномерной загрузки системы определяется методом их выбора.

В рассматриваемом случае модель массового обслуживания зависит от четырех параметров: λ ; μ_1, μ_2 — интенсивностей обслуживания каналов; φ — вероятности выбора более «быстрого» или более «медленного» канала. Эффективность системы при данном значении φ зависит только от двух параметров:

$$\psi = \frac{\lambda}{\mu} = \frac{\lambda}{\mu_1 + \mu_2}, \quad (4)$$

и отношение интенсивностей

$$\alpha = \frac{\mu_2}{\mu_1}, \quad (5)$$

величина которого меньше 1.

На основании графа (рис. 4) состояний составим матрицу вероятностей перехода (4) и систему дифференциальных уравнений. Прежде всего рассмотрим возможные состояния системы и их вероятности. Необходимо отличать случай, когда занят прибор 1, от случая, когда занят прибор 2. Обозначим соответствующие вероятности через p_1, p_2 . Так как при $n \neq 1$ обозначение p_n не приводит к недоразумению, мы сохраним его. Соответственно, возможные вероятности состояния системы обозначаются как

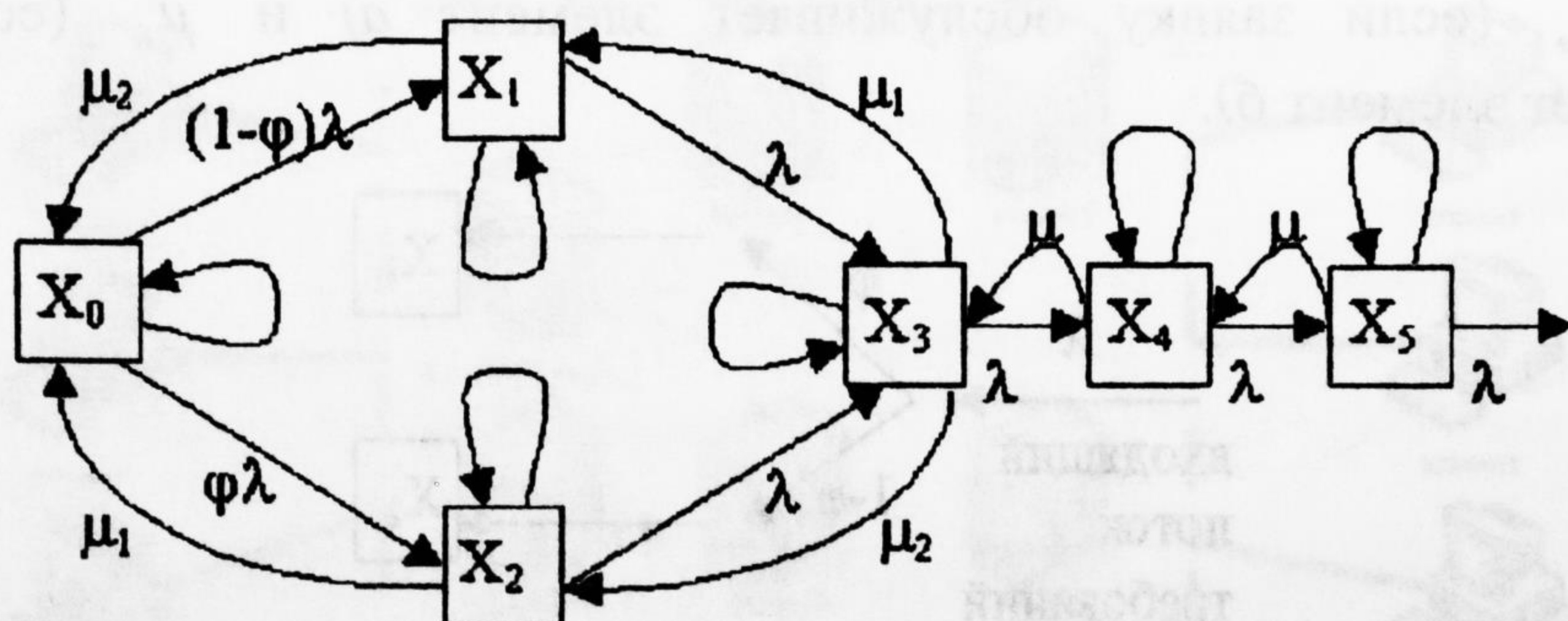


Рис. 4. Граф состояний

Следует отметить, что состояний может быть больше, чем представлено на графе (рис. 4), однако нас интересуют только состояния, связанные с приборами X_0, X_1, X_2, X_3 .

Уравнения состояния имеют следующий вид:

$$p_0(t+dt) = (1-\lambda dt)p_0(t) + \mu_1 p_1(t)dt + \mu_2 p_2(t)dt, \quad (6)$$

$$p_1(t+dt) = [1-(\lambda + \mu_2)dt] p_1(t) + \mu_1 p_3(t)dt + (1-\varphi)\lambda p_0(t)dt, \quad (7)$$

$$p_2(t+dt) = [1-(\lambda + \mu_1)dt] p_2(t) + \mu_2 p_3(t)dt + \varphi\lambda p_0(t)dt, \quad (8)$$

$$p_n(t+dt) = [1-(\lambda + \mu)dt] p_n(t) + \mu p_{n+1}(t)dt + \lambda p_{n-1}(t)dt, n \geq 2 \quad (9)$$

Таблица 1

Матрица состояний в момент $t+dt$

		Состояние в момент $t+dt$					
		n=0	n=1 (x1)	n=1 (x2)	n=2	n=3	
Состояние в момент t	n=0	$1-\lambda dt$	$(1-\varphi)\lambda dt$	$\varphi\lambda dt$	0	0	...
	n=1 (x1)	$\mu_2 dt$	$1-(\lambda + \mu_2)dt$	0	λdt	0	...
	n=1 (x2)	$\mu_1 dt$	$1-(\lambda + \mu_1)dt$	0	λdt	0	...
	n=2	0	$\mu_1 dt$	$\mu_2 dt$	$1-(\lambda + \mu)dt$	λdt	...
	n=3	0	0	0	μdt	$1-(\lambda + \mu)dt$...

Цепь, соответствующая матрице (табл. 1), является неприводимой и неперiodической. В [4] доказано, что в таком случае все вероятности $p_n(t)$ имеют пределы p_n и, кроме того, эти пределы p_n либо все положительны, либо

все равны нулю. При $t \rightarrow \infty$ производные, входящие в уравнения состояния (6-9), стремятся к нулю. Отсюда уравнения (6-9) в установившемся режиме принимают вид:

$$\lambda p_0 = \mu_2 p_1 + \mu_1 p_2;$$

$$(\lambda + \mu_2) p_1 = \mu_1 p_3 + (1 - \varphi) \lambda p_0;$$

$$(\lambda + \mu_1) p_2 = \mu_2 p_3 + \varphi \lambda p_0;$$

$$(\lambda + \mu) p_n = \mu p_{n+1} + \lambda p_{n-1}, n \geq 2.$$

Введя отношение интенсивностей обслуживания α , (4,5), получим

$$p_1 = \frac{\psi}{1+2\psi} \cdot \frac{1+\alpha}{\alpha} (\psi+1-\varphi) p_0, \quad (10)$$

$$p_2 = \frac{\psi}{1+2\psi} \cdot (1+\alpha)(\psi+\varphi) p_0, \quad (11)$$

$$p_3 = \frac{\psi^2}{1+2\psi} \cdot \frac{1+\alpha}{\alpha} [1+(1+\alpha)\psi - (1-\alpha)\varphi] p_0, \quad (12)$$

$$p_n = \frac{\psi^n}{1+2\psi} \cdot \frac{1+\alpha}{\alpha} [1+(1+\alpha)\psi - (1-\alpha)\varphi] p_0 \quad (13)$$

Учитывая то, что сумма вероятностей всех состояний равна 1

$$p_0 = \frac{1-\psi}{1 + \frac{\psi}{1+2\psi} \cdot \frac{1}{\alpha} [1+(1+\alpha^2)\psi - (1-\alpha^2)\varphi]} \quad (14)$$

Таким образом, из формул (10-14), можно определить вероятности отклонения системой запроса клиента для каждой пары узлов.

Исходя из вышесказанного, определим алгоритм выбора узла системы балансером для привязывания клиента по его запросу (рис. 5).

1. При регистрации сервера в балансере сохраняются его статические параметры (3).

2. Через каждый временной интервал (t) происходит опрос всех зарегистрированных в балансере узлов (рис. 5). Во время опроса с узла поступают динамические параметры (3). Полученные параметры сохраняются во временной памяти для последующего использования.

3. При поступлении клиента в систему балансер производит попарное сравнение узлов системы. Для каждого узла из временного хранилища определяется интенсивность обслуживания μ , вероятности φ и общее количество обслуживаемых клиентов $n = uCount_i + uCount_{i+1}$.

4. Из формул (10-14) вычисляем вероятность отказа в обслуживании для этих двух узлов.

5. Выбирается пара, где вероятность отказа минимальная, далее случайным образом назначается один из узлов пары (вероятность выбора узла задается параметром φ).

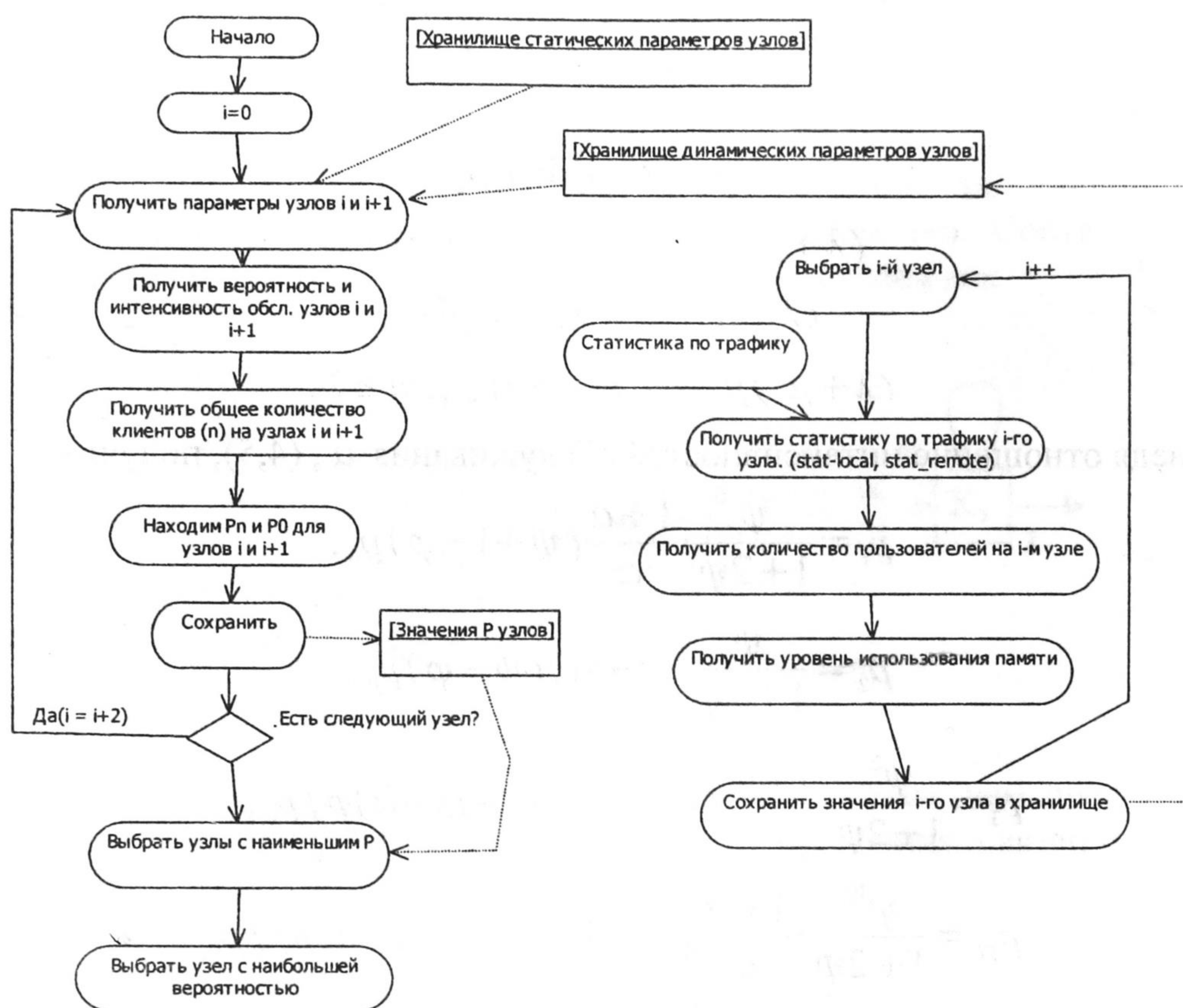


Рис. 5. Алгоритм определения «оптимального» узла

Вывод

Предложен алгоритм работы балансера в кластерной системе реализующие вебинары, основанный на теории массового обслуживания.

СПИСОК ЛИТЕРАТУРЫ

1. Википедия. Свободная энциклопедия // <http://ru.wikipedia.org/wiki/>.
2. Мостицкий И. Современные английские термины из области электроники. Вып. 34 (224) // Электронная энциклопедия. 2004.
3. Вентцель Е.С. Теория вероятностей. М.: Государственное издательство физико-математической литературы, 1962. 564 с.
4. Кофман А., Крюон Р. Массовое обслуживание. Теория и приложения. М.: Мир, 1965. 302 с.
5. Олзоева С.И. Моделирование и расчет распределенных информационных систем. Улан-Удэ: Изд-во ВСГТУ, 2004. 67 с.