

# Основы теории информации и кодирования

14 июля 2024 г.

## Приложение А. Регламент курса

В данном разделе даны базовые положения регламента курса ОТИК. **Все особые случаи обсуждаются с преподавателем индивидуально**, решение в каждом случае принимается отдельно и только **после воплощения соответствующего случая в реальность**. То, что было возможно в начале семестра — недопустимо в конце.

Хотя декабрь загружен не так сильно, как май, но всё же гораздо сильнее октября.

### А.1. Итоговая оценка $R$

Оценка («неуд»/«удовл»/«хор»/«отл») выставляется по совокупности:

- выполненной в семестре работы;
- обязательного для всех экзамена (или зачёта у НБ-3\*).

Баллы в ОРИОКС подгоняются под полученную оценку за счёт граф «Итог» и «Натяжка»; при этом **честно заработанные в семестре баллы не уменьшаются**. Баллы, полученные жульничеством (плагиат, манипуляции с размером команды и т. д. — полный список не приводится, так как нельзя объять необъятное), **удаляются** как техническая ошибка — независимо от того, сколько времени они пробыли в ОРИОКС.

#### А.1.1. ОРИОКС: сумма баллов $S$ и формирование итоговой оценки $R$

Текущая оценка  $R$  записывается в ОРИОКС в виде суммы баллов  $S$ .

Непосредственно на баллы  $s$  (линейный рост  $S_{\text{новое}} = S_{\text{старое}} + s$ ) оцениваются:

- лабораторные работы, защищённые на 1–16 неделях по графику (раздел **РЛ**);
- задачи, решаемые на лекциях (раздел **А.2.1**);
- работа на семинарах (раздел **А.2.1**);
- замечания и дополнения (раздел **А.6**);

На оценку  $r$  (нелинейный расчёт ( $R_{\text{старое}}$  и  $r$ )  $\rightarrow R_{\text{новое}} \rightarrow S_{\text{новое}}$ ) оцениваются:

- все альтернативные лабораторным работам варианты (раздел **А.5**);
- экзамен/зачёт (раздел **А.3**).

Не принимаются лабораторные работы, сдаваемые: помногу сразу, на 17–18 неделях или в сессию (подробнее — в приложении **Б**, раздел **РЛ**).

#### Шкала оценок

Внутри семестра (как для итоговой  $R$ , так и для КМ  $r$ ) как особый случай рассматривается оценка «3 с натяжкой» («неуд» < «3 с натяжкой» < «удовл»). В итоговой ведомости «3 с натяжкой» отображается как «удовлетворительно».

Сумма баллов  $S$  и оценка  $R$  связаны соотношением:

$$R(S) = \begin{cases} \text{«отл»}, & S \geq 86, \\ \text{«хор»}, & 70 \leq S < 86, \\ \text{«удовл»}, & 50 \leq S < 70 \text{ и графа «Натяжка» пуста/0/н,} \\ \text{«3 с натяжкой»}, & S \approx 50 \text{ и графа «Натяжка» не пуста,} \\ \text{«неуд»}, & S < 50. \end{cases} \quad (\text{A.1})$$

Оценка «3 с натяжкой» не может быть повышена до «хор».

### А.1.2. Графы КМ в ОРИОКС

Так как для получения оценки по ОТИК есть несколько альтернативных вариантов, иногда взаимоисключающих — многие из граф КМ в ОРИОКС могут быть пустыми даже у студента с итоговой  $R = \text{«отл»}$ . Графа «Натяжка» у нормально успевающего студента *должна* быть пустой.

#### Баллы и специальные значения граф

В графу КМ ОРИОКС могут быть выставлены значения:

- а)  $s \geq 1$ : КМ успешно защищено с баллами  $s$ ;
- б)  $s = 0$  или  $s = \text{н}$ : защиты обязательного КМ не было;
- в)  $s = 0,01$ : маркер «защита запрещена» — выставляется, если:
  - сдан альтернативный вариант, несовместимый с КМ (раздел А.5);
  - защита КМ уже была: сдавали на повышение итоговой оценки  $R(S)$ , но в итоге подтвердили старую;
- г)  $s = 0,02$ : маркер «защита провалена» — оценка за КМ «неуд».

КМ с 0/н можно сдать любому преподавателю ОТИК, кто согласится принимать, и повысить оценку; 0,02 — либо тому, кто выставил 0,02, либо комиссии; 0,01 — обсуждается индивидуально и только с тем, кто выставил.

#### Графы лабораторных работ и КР — значение выставляется один раз

В графы лабораторных работ  $L_i$ , Кр1 и Кр2 значение  $s$  выставляется после защиты (при повышении — заменяется на новое).

Для работы  $L_i$ , защищённой на 1–16 неделях по графику:

- а) базовые баллы  $s \geq 1$  выставляются в графу « $L_i$ »;
- б) бонусные баллы  $s^{\text{bonus}} \geq 0$  добавляются в графу «Бонус (л/р)».

#### Накопительные графы бонусов — значение растёт от 0 в течение семестра

Со значениями граф «Зд (лек)», «Бонус (прочее)», «Бонус (л/р)» и т. п. вновь заработанные баллы суммируются ( $s_{\text{накоп. новое}} = s_{\text{накоп. старое}} + s$ ):

- а) зд (лек) — задачи, решаемые на лекциях, выставляет лектор;
- б) бонус (прочее) — иные баллы, которые назначает лектор и для которых нет отдельной графы (в частности, за замеченные опечатки и ошибки);
- в) бонус (л/р):
  - бонусные задания всех лабораторных работ  $L_i$ ;

- иные бонусные баллы, которыми преподаватель отмечает какую-либо  $L_i$ ;
- а также все прочие баллы, которые назначает преподаватель лабораторных работ (на своё усмотрение) и для которых нет отдельной графы.

### Графы семинаров — два раза в семестр

Графа «Зд (сем) <недели>» — баллы за работу на семинарах на соответствующих неделях, без учёта посещаемости; графа «П (сем) <недели>» — отдельно посещаемость семинаров. Значения  $s$  выставляются после окончания интервала:

- а) за 1–8 недели — на 9 неделе,  $0 \leq s \leq \max(s)$ , где  $\max(s)$  — ёмкость графы;
- б) за 9–16 недели — дважды, предварительные (9–14) и полные (9–16) значения:
  - на 15 неделе — предварительные  $0 \leq s \leq \frac{3}{4} \max(s)$ ;
  - на 17 неделе — полные  $0 \leq s \leq \max(s)$ .

Баллы за каждый интервал (1–8 или 9–16 недели) выставляются независимо друг от друга. Оценивается интервал в целом, «баллы за один плюсики» не определены.

К баллам посещаемости: если студентов мало и я могу опросить всех, то «пришёл просто так посидеть» оценивается как «не был».

### Графы «Итог», «Курсовая» и «Индивидуальное задание» — значение графы подгоняет общую сумму баллов $S$ под оценку $R_{\text{Итог}}$ (нелинейный расчёт)

Если при текущей сумме баллов  $S_{\text{старое}}$  оценка на экзамене  $R_{\text{Итог}} \geq R(S_{\text{старое}})$ , то значение  $s$  графы «Итог» задаётся так, что  $R(S_{\text{старое}} + s) = R_{\text{Итог}}$ . Если ёмкости графы «Итог» не хватает — используется также «Бонус (прочее)».

Если  $R_{\text{Итог}} < R(S_{\text{старое}})$  — в зависимости от ситуации возможно как сохранение текущих баллов  $S_{\text{старое}}$  и  $s = 0,02$  в «Итог», так и поиск ошибки в ОРИОКС (для чего студенту необходимо проявить не только исключительные знания, но также настойчивость и целеустремлённость: такие случаи исчезающе редки).

Аналогично рассчитываются баллы граф «Курсовая» или «Индивидуальное задание» (если их нет — используется графа «Бонус (прочее)»).

### Графа «Натяжка»

Если оценка  $R_{\text{Итог}} = \langle 3 \text{ с натяжкой} \rangle$ , то дополняющие до 50 баллы выставляются в графу «Натяжка»: они не соответствуют никаким реальным достижениям студента. Если студент повышает  $R_{\text{Итог}}$  до «удовл» и выше — выставляются баллы за сделанное, а графа «Натяжка» обнуляется.

## А.2. Работа в семестре (1–16 недели)

Напомню тот очевидный факт, что работа в семестре включает не только *контроль* знаний, умений и навыков, но и их *получение*. Это невозможно без обратной связи, поэтому, если что-то непонятно — **спрашивайте!!!** Даже перед экзаменом всегда есть консультация для разрешения вопросов по курсу.

Студент, приступающий к изучению курса ОТИК, по определению ещё его не знает (знающие — могут доказать это в начале семестра и более не появляться, см. раздел А.5.3).

### Нормальное изучение курса — поэтапное

Лабораторная работа  $L_i$  (приложение Б, регламент — в разделе РЛ) выполняется сразу же после того, как её тема рассказана на лекции. Возникшие вопросы разрешаются на лабораторном занятии, на семинаре или после следующей лекции; после их разрешения следует защита  $L_i$ .

Баллы (базовые и бонусные) за лабораторные работы характеризуют *процесс* изучения ОТИК в течение семестра (с 1 по 16 неделю). Баллы за семинары также характеризуют процесс изучения, а не только правильность ответов.

Если оценка  $R(S)$  по ним не соответствует знаниям студента — можно её изменить, оценив *результат* (раздел А.5).

### Авральное изучение не делится на этапы

Если студент приступил к изучению курса ОТИК на 13 неделе или позже, процесс изучения невозможно разделить на сколько-нибудь длительные этапы. Сдавать лабораторные работы в такой ситуации **нельзя**.

В этом случае оценивается *результат* изучения по одному из альтернативных вариантов (раздел А.5). Если при подготовке возникли сложности — вопросы также необходимо разрешить на лабораторном занятии, на семинаре или после лекции.

#### А.2.1. Занятия в течение семестра

Практика показывает, что если с кодами по умолчанию, изучаемыми в курсе ОЭВМ [и Асм], студенты могут более или менее успешно разобраться самостоятельно — то с кодами и методами сжатия ОТИК у тех, кто не ходил ни на лекции, ни на семинары, возникают большие сложности.

Статьи, как правило, иллюстрируют метод на примере *одного* реализующего его алгоритма и одного кода, а в вариантах лабораторных/контрольных работ есть *разные* коды и соответствующие им алгоритмы, реализующие один и тот же метод.

#### Лекции

Посещаемость поточных лекций не оценивается.

Материалы лекций выложены в <https://gitlab.com/illinc/otik> — этого года — в корне, прошлых лет — в соответствующих поддиректориях.

#### Семинары

Посещаемость семинара оценивается только если занятие действительно проведено как семинар, а не выродилось в индивидуальную консультацию.

Задачи для семинаров для подготовки к ним вписаны в презентации лекций. Решения с уже прошедших семинаров верстаются и выкладываются только если попросит тот, кто решал, причём *до* стирания с доски (файл [sem44-rs\\_2022.pdf](#) — пока единственный такой случай).

Пропущенные даже по уважительной причине семинары «отработать» нельзя; можно придти с другой группой. Вопросы по пропущенным темам задавайте на консультациях, баллы добивайте альтернативными вариантами (раздел [А.5](#)).

### Лабораторные занятия

На ближайшем занятии после самостоятельной работы над  $Li$  необходимо:

- разрешить с преподавателем все вопросы, которые возникли при подготовке;
- защитить выполненные задания  $Li$ .

### Консультации

На 1–16 неделях каждый преподаватель регулярно проводит общие консультации. Расписание и надо ли записываться заранее — уточняйте у преподавателя.

На консультациях можно как задавать любые вопросы по курсу ОТИК, так и сдавать лабораторные работы, но приоритет у обсуждения ВКР/диссертаций выпускных курсов.

## А.3. Экзамен/зачёт

Экзамен обязателен для всех студентов. Зачёт у НБ-3\* проводится аналогично экзамену у ПИН-3\* и тоже обязателен, поэтому не описывается отдельно.

К экзамену/зачёту **по расписанию допускаются все** (в том числе возможно повышение оценки с нуля до «3 с натяжкой» — написанием Кр1 и Кр2 и ответами). В этот день разрешаются все сложные вопросы. О досрочном см. раздел [А.3.2](#).

Ответ оценивается не на баллы, а на оценку «неуд»  $\leq R_{\text{итог}} \leq$  «отл», баллы добавляются до необходимого значения  $S$  (см. раздел [А.1.2](#)).

**Зачтение статей из Интернета и даже текста лекций ответом на вопрос не считается; оценивается как «неуд».** Говорите своими словами. Поясняйте сказанное на примерах.

### А.3.1. Билеты и базовые вопросы

Списка базовых вопросов нет и не будет, так как их цель — проверить вашу логику, а не способность к заучиванию. Список билетов выложен в <https://gitlab.com/illinc/otik> и мало меняется с годами.

### А.3.2. Досрочная сдача экзамена/зачёта

Студенты, допущенные к досрочному экзамену/зачёту (таблица А.1), отвечают только на базовые вопросы. Для допуска к досрочному — необходимо выполнить лабораторные работы **по графику** в указанном объёме.

Также к досрочному экзамену/зачёту допускаются студенты, успешно выполнившие курсовую работу (раздел А.5.3) или индивидуальное задание (А.5.2).

#### Допуск к досрочному экзамену/зачёту

Таблица А.1

	Недели	Дата	Допуск
<b>Экзамен ПИН-3* (180/5) — ОТИК</b>			
Досрочный-1	15–16	в день последней лекции	Л1–Л5 и зачётка
Досрочный-2	17–18	назначается лектором	Л1–Л4 и зачётка
Плановый	сессия	по расписанию экзаменов	—
<b>Зачёт НБ-31 (108/3) — ОТИК</b>			
Досрочный	15–16	назначается преподавателем	Л1–Л2 и зачётка
Плановый	17–18	назначается преподавателем	—

Для НБ оценка осенью не выставляется в ведомость ОТИК, но записывается и учитывается весной в ведомости ОЭВМ (ОТИК НБ включает часть тем ОЭВМ).

### А.4. Долги и пересдачи

Итоговую оценку  $R_{\text{Итог}}$  = «неуд» можно повысить до конца сессии — новая оценка попадает в дополнительную ведомость; иначе образуется **долг**.

Закрыть долг можно **до конца следующего семестра**. Для этого **заранее, до** официальной даты пересдачи, указанной в ОРИОКС, необходимо найти преподавателя, получить и выполнить задание.

Начиная со второго после экзамена семестра (следующая весна и позже) долг не может быть закрыт без письменного разрешения ответственного кафедры.

Для закрытия долга необходимо:

- либо (на «неуд»  $\leq R_{\text{Итог}} \leq$  «удовл») написать Кр1 и Кр2 (раздел А.5.1);
- либо (на «неуд»  $\leq R_{\text{Итог}} \leq$  «отл») получить и выполнить индивидуальное задание (раздел А.5.2).

Должник не имеет права сдавать тот же набор лабораторных работ, что и его однокурсники ранее, если только преподаватель не даст ему такое задание.

#### А.4.1. Долг ОЭВМ [и Асм] с весны

Для закрытия долга ОЭВМ [и Асм] *достаточно* защитить **ЛО**, написать Кр0 (КрЭВМ) и ответить на базовые вопросы ОЭВМ (все три пункта). Ответить на базовые вопросы и выполнить хотя бы одно из **ЛО/Кр0** — *необходимо*, но может оказаться недостаточным в зависимости от ситуации: ответов по ОЭВМ, качества защиты **ЛО/Кр0**, успехов на семинарах ОТИК.

### А.5. Альтернативы лабораторным работам

По результатам защиты каждого из вариантов выставляется **итоговая оценка за весь курс** «неуд»  $\leq R_{\text{Итог}} \leq R_{\text{max}}$  (баллы  $S$  подгоняются под оценку  $R_{\text{Итог}}$ , раздел **А.1.2**; все обязательные графы «Л $i$ » маркируются  $s = 0,01$ ).

В исключительных случаях, если исполнение и защита безукоризнены, а причина опоздания уважительна — допускается  $R_{\text{Итог}} = R_{\text{max}} + 1$ : «хор» для Кр1 и Кр2 и «хор»/«отл» для индивидуального задания на «удовл»/«хор».

#### А.5.1. Кр1 (КрХф) и Кр2 (КрRLE/LZ)

$R_{\text{max}} = \text{«удовл»}$  за обе совокупно. По отдельности — оцениваются на баллы соответственно качеству КР и ёмкости графы.

Задания выложены в <https://gitlab.com/illinc/otik>; вариант назначается преподавателем (как из выложенных в <https://gitlab.com/illinc/otik>, так и индивидуальный). Можно писать командой. Ответ без решения = «неуд».

Должникам ПИН по ОЭВМ [и Асм] необходимо написать все три: КрЭВМ (Кр0), Кр1 и Кр2. Студентам НБ — КрЭВМ (Кр0) и Кр1.

#### А.5.2. Индивидуальные задания

$R_{\text{max}} \in \{\text{«удовл»}, \text{«хор»}\}$ , в зависимости от сложности (графа «Индивидуальное задание», если её нет — «Бонус (прочее)»). Объём задания должен соответствовать оставшемуся от семестра времени; на 17 неделе и позже — не более объёма двух лабораторных работ.

#### Блиц-задания

$R_{\text{max}} \in \{\text{«удовл»}, \text{«хор»}\}$ , в зависимости от сложности.

Сложность блиц-заданий на порядок меньше сложности тех, что даются на остаток семестра, но выполнить их требуется немедленно, не выходя из аудитории.

#### А.5.3. Курсовая работа (командная или индивидуальная)

Так как мне надоело получать в конце сессии нечто, что не тянет даже на лабораторную работу, но предлагается как «курсовая» — этот путь для ОТИК



закрыт. Желающие меня переубедить — подходите лично, с планом конкретной работы.

## А.6. Замечания и дополнения

Замечания и дополнения к данному документу можно отправить в письменном виде по адресу <https://gitlab.com/illinc/otik/issues>.

Принятое замечание/дополнение приносит первому приславшему его студенту от 1 до 8 бонусных баллов.

## А.7. Обновление пособия

Между семестрами задания и регламент лабораторных работ качественно обновляются. Выполнять лабораторные работы прошлого года — **нельзя**.

В течение семестра:

- уточняются **формулировки** заданий и регламента, если исходные вызывают недопонимание; суть заданий и регламента — неизменна до конца семестра;
- появляются **новые бонусные** задания;
- дополняется теория (в том числе по вашим замечаниям и дополнениям); поэтому периодически скачивайте актуальную версию пособия из репозитория <https://gitlab.com/illinc/otik>. Лабораторные работы в ОРИОКС — неактуальны по определению из-за сложной процедуры обновления. Прямая ссылка на актуальную версию: [https://gitlab.com/illinc/otik/-/raw/master/otik-labs.pdf?ref\\_type=heads&inline=false](https://gitlab.com/illinc/otik/-/raw/master/otik-labs.pdf?ref_type=heads&inline=false)

## А.8. Примечание для студентов, изучающих по индивидуальному плану ОТИК раньше, чем ОЭВМ [и Асм], что противоречит рабочей программе

Прежде чем изучать методы сжатия и защиты от помех (ОТИК), необходимо подтвердить знакомство с кодами, которые используются в ЭВМ по умолчанию (ОЭВМ [и Асм], <https://gitlab.com/illinc/gnu-asm>).

Тем, кто ещё не имеет оценки за ОЭВМ [и Асм], необходимо:

- перед Кр1 и Кр2 по ОТИК написать КрЭВМ (Кр0) из ОЭВМ [и Асм];
  - перед защитой Л2 защитить Л0 (Л1 может защищаться параллельно с Л0).
- Для должников ОЭВМ [и Асм] Л0 + Кр0 + базовый вопрос = закрытие долга.

Так как у групп НБ-3\* ошибочное расположение (ОТИК раньше ОЭВМ) не будет исправлено в ближайший год, на лекциях и семинарах ОТИК вначале разбираются коды по умолчанию.

## Приложение Б. Лабораторный практикум по основам теории информации и кодирования

Регламент курса в целом описан в приложении А.

### РЛ. Регламент лабораторных работ

Баллы (базовые и бонусные) за лабораторные работы  $L_i$  характеризуют *процесс* изучения ОТИК в течение семестра, поэтому их можно получить, только сдавая  $L_i$  **по графику**:

- начиная с первого или второго лабораторного занятия (Л1 и Л2 могут выполняться даже до первой лекции, они не требуют специфических знаний);
- по одной работе на каждом занятии, максимум — две на одном занятии;
- с минимальным опозданием.

По результатам защиты лабораторной работы  $L_i$  выставляются:

а) на 1–16 неделях по графику (раздел РЛ.4):

- базовые баллы  $0 \leq s \leq s_{\max}$  в графу « $L_i$ »;
- бонусные баллы  $0 \leq s^{\text{bonus}} \leq s_{\max}^{\text{bonus}}$  в «Бонус (л/р)» (раздел РЛ.3);

б) на 1–16 неделях в составе комплекта из трёх и более работ (раздел РЛ.6):

- баллы  $0 \leq s \leq 1$  в графе « $L_i$ »;
- нет бонусных баллов:  $s^{\text{bonus}} = 0$ ;

в) 17–18 недели, сессия — лабораторные работы не принимаются (раздел РЛ.7).

Базовые  $s$  и бонусные  $s^{\text{bonus}}$  баллы выставляются преподавателем:

- не за программу или отчёт, а за **защиту** лабораторной работы;
  - на его усмотрение: если описанных явно бонусов/штрафов недостаточно для адекватного оценивания конкретной защиты — вводятся дополнительные;
  - сообразно **качеству** защиты и работы:  $s$  при  $0 \leq s \leq s_{\max}$  может быть и 1, и 0.
- Из баллов  $L_i$ , просроченной более чем на две недели без уважительной причины, вычитается величина опоздания (таблица РЛ.1). Просроченная на  $\Delta t \geq 12$  недель даже по уважительной причине — не может быть сдана вообще.

#### РЛ.1. Командная работа

Состав команды: либо один студент, либо 2–3 сидящих рядом студента. Задания для  $n \in \{1, 2\}$  одинаковы; для  $n \geq 3$  есть дополнительные (для троек). Независимо от количества участников  $n$ , команда выполняет **один вариант** работы — по номеру команды (№), один отчёт и **одну защиту**, получает **одну оценку**.

Каждый из соавторов должен уметь объяснить все результаты лабораторной работы и модифицировать свою часть кода (таблица РЛ.2). Студента, который

## График снижения базовых баллов за лабораторные работы

Таблица РЛ.1

ОТИК: группы ПИН-31–32								
Неделя	max Л1	max Л2	max Л3	max Л4 (упрощ.)	max Л5 (упрощ.)	max Л6 (упрощ.)	max 3–4 шт. сразу	max ≥ 5 шт. сразу
1, 2	7						3	1
3, 4	7	7					3	1
5, 6	6	7	7				3	1
7, 8	5	6	7	14 (7)			3	1
9, 10	4	5	6	14 (7)			3	1
11, 12	—	4	5	13 (6)	14 (7)		3	1
13, 14	—	—	4	12 (5)	14 (7)		1	1
15, 16	—	—	—	11 (4)	13 (6)	14 (7)	1	1
15, 16 (с нуля) 17, 18 сессия	л/р не принимаются — пишите Кр1+Кр2 или выполняйте индивидуальное задание							

Максимумы граф включают также некоторые бонусы:

№	Л1	Л2	Л3	Л4	Л5	Л6
max базовой части	7	7	7	14	14	14
max бонуса внутри графы	0	+7 Маркова	+7 папки	+14 АС	+14 LZ77	нет графы (бонусная л/р)
max графы	7	14	14	28	28	

## Первый семестр ОЭВМ/ОТИК (в ОРИОКС — ОТИК): группы НБ-3\*

Неделя	max Л1	max Л0	max Л2	max Л3	max 3 шт. сразу	max ≥ 4 шт. сразу
1, 2, 3, 4	16				3	1
5, 6, 7, 8	16	16			3	1
9, 10, 11, 12	—	16	16		3	1
13, 14, 15, 16	—	—	16	16	1	1
13–16 (с нуля) 17, 18 сессия	л/р не принимаются — пишите Кр0+Кр1 или выполняйте индивидуальное задание					

не только не пришёл на защиту  $L_i$ , но и не помогал её делать, присутствующие соавторы могут временно (на  $L_i$ ) или постоянно исключить из команды.

Команды из  $n \geq 4$  студентов запрещены на ВЦ МИЭТ; в не-ВЦ классах (и задания для  $n \geq 4$  сверх заданий для троек) — на усмотрение преподавателя.

Командная работа и сложные ситуации

Таблица РЛ.2

	Нормально	Недопустимо
<b>Изменение состава команды</b> (постоянное)	Один-два раза (не сработались) — баллы за сданные $Li$ сохраняются. № команд д. б. уникальным — утвердите № у преподавателя, прежде чем делать новый вариант	Постоянно делиться/объединяться синхронно с наличием/отсутствием заданий для троек — штраф ко всем $Li$ по $-2$ балла к каждой
<b>Защита <math>Li</math> в неполном составе</b> (подготовка $Li$ — всегда в полном)	Отсутствовать на защите некоторых $Li$ — защищаться отдельно не нужно. Тройка выполняет задания для троек, даже если на защите один студент	Отсутствовать систематически (и не сообщать о сложностях): у пропавшего на полсеместра без предупреждения — аннулируются баллы за все сданные без него $Li$
<b>Исключение из команды на одну <math>Li</math></b> (временное)	Делать и защищать с командой $L(i + 1)$ и следующие. На тему $Li$ — получить и выполнить доп. задание или смириться с $s(Li) = 0$	Защищать ту же самую $Li$ , которую уже защитила команда — в лучшем случае будет $s(Li) = 1$
Постоянное исключение из команды — изменение состава; исключённый меняет №		
<b>Разные мнения на защите <math>Li</math></b>	Прийти к согласию, затем отвечать. Не удалось договориться — скажите об этом и изложите все мнения последовательно	Излагать все мнения одновременно — ни один ответ не засчитывается; озлобленность преподавателя растёт
<b>Разные мнения при подготовке <math>Li</math></b>	Спросить у преподавателя и только потом защищаться	Всё остальное — штрафуются либо как недоделка, либо как жульничество

РЛ.2. Требования к выполнению лабораторных работ

Лабораторные работы могут быть подготовлены дома и защищаться на ВЦ МИЭТ. Перед защитой на лабораторном занятии необходимо разрешить все вопросы, возникшие в процессе подготовки.

Средства разработки

Любые, позволяющие работать с «сырыми» бинарными данными, в том числе заданной разрядности. Если вы записываете в файл два байта — шестнадцатеричный просмотрщик/редактор должен показать, что там именно два байта, причём именно те, что вы записывали.

Если в желаемом языке есть нужные средства, но вы ими не владеете — либо изучайте их сейчас, либо пишите на C/C++, коллекция компиляторов GCC.

## Отчёт

Дистанционные группы оформляют отчёт в любом случае. Для очных групп:

- если  $L_i$  сделана на занятии — защищайте её сразу, **не оформляя отчёта**;
- если  $L_i$  готовится дома заранее — параллельно оформляйте отчёт.

Формат — OpenDocument или PDF. Заголовок отчёта должен включать имя группы и ФИО авторов, тему работы, а также — для каждого задания:

- номер и текст задания;
- номер и текст варианта (если есть);
- **платформа**, на которой делается задание, в частности — ОС и разрядность;
- ключевые фрагменты программного кода;
- ключевые пояснения, которые вы боитесь забыть;
- ссылки на использованные справочные материалы и пояснения, что и для чего вы там искали, где применили найденное;
- результат выполнения задания: результаты измерений с комментариями.

Вместе с отчётом должен предоставляться полный текст программ, готовый к сборке и запуску, поэтому копировать их ещё и в отчёт не нужно.

## РЛ.3. Необязательные (бонусные) задания (1–16 недели)

Задания, отмеченные как **«Бонус»**, *необязательны*. На 1–16 неделях за их выполнение начисляются дополнительные баллы  $0 \leq s^{\text{bonus}} \leq s_{\text{max}}^{\text{bonus}}$  (максимальное количество  $s_{\text{max}}^{\text{bonus}}$  указано в тексте задания и не уменьшается со временем:  $s^{\text{bonus}}$  зависит только от качества исполнения и защиты, но не от времени сдачи).

Баллы  $s^{\text{bonus}}$  добавляются к графе «Бонус (л/р)». Бонусные задания могут быть сданы либо одновременно с соответствующей лабораторной работой, либо (при условии, что это не помешает преподавателю принимать у других студентов обязательные задания) после неё.

При защите нескольких лабораторных работ сразу (раздел **РЛ.6**) бонусные баллы не начисляются ( $s^{\text{bonus}} = 0$ ), бонусные задания формируют базовые баллы  $s$  наравне с обязательными.

## РЛ.4. Защита лабораторной работы (1–16 недели)

До начала защиты лабораторной работы  $L_i$  команда должна задать преподавателю подготовленные во время выполнения  $L_i$  вопросы.

Если у вас вопрос — говорите «у меня вопрос», если планируете защищать — «хочу сдать»/«хочу защитить работу». Запрос «посмотрите лабораторную работу» будет по возможности игнорироваться, так как его цель не ясна.

Дополнительные вопросы на защите задаются всем. **Отчёта для оценивания недостаточно.** По итогам защиты ставятся баллы  $s + s^{\text{bonus}}$ .

### Защищать не по порядку — можно, $\max(s_{\max})$ корректируются

Если какую-то  $L_i$  не удалось выполнить даже частично:

- сформулируйте вопросы к преподавателю о том, что вызвало проблемы в  $L_i$ ;
- если  $L(i+1)$  сделать получается — делайте и защищайте;
- на том же занятии задавайте вопросы про  $L_i$ ;

$\max(s_{\max})$  считается по фактическому порядку защит, а не по номерам  $i/(i+1)$ .

### Не выполненная, но проработанная $L_i$ + дополнительное задание = «удовл»

Если ни одно задание лабораторной работы  $L_i$  не выполнено, но:

- по всем обязательным заданиям  $L_i$  подготовлены вопросы к преподавателю;
- сохранены все неудачные попытки выполнения заданий с комментариями, чем именно они плохи (не собирается — с какими сообщениями компилятора?

Не приводит к нужному результату — а к какому приводит?);

- после разрешения вопросов команда выполнила по указанию преподавателя либо задание  $L_i$ , либо дополнительное задание, не покидая аудиторию;

то  $L_i$  засчитывается на  $1 \leq s \leq \frac{1}{2} \max(s_{\max})$ , без бонусных баллов ( $s^{\text{bonus}} = 0$ ).

### Чего делать не стоит («неуд», то есть 0 баллов)

Однозначно на «неуд» (на  $s = s^{\text{bonus}} = 0$  баллов) оценивается работа  $L_i$ , если вопросов перед защитой не задавалось, зато на самой защите:

- звучат ответы «делал давно, не помню», «взял в Интернете, не помню где»;
- есть хотя бы одна программа (в том числе отлично выполненная), в которой студенты ничего не могут объяснить и/или изменить;
- в отчёте есть хотя бы одно место, которое студенты не могут разъяснить и пересказать своими словами;
- хотя бы в одном задании  $L_i$  (основном или бонусном) есть признаки плагиата (несоответствие варианту, несоответствие программы заявленным ОС и разрядности, использование заданий предыдущего года и ранее, и т. п.);

как только хоть что-то из этого проявилось — защита завершается. Преподаватель может как дать команде дополнительное задание на  $1 \leq s \leq 2$  (но  $s^{\text{bonus}} = 0$ ), так и сразу отправить на пересдачу (на следующем занятии, тому же преподавателю).

### РЛ.5. Дистанционная защита для дистанционных групп

Дистанционные группы защищают лабораторные работы, используя:

- синхронную связь с демонстрацией экрана (телемост=видеоконференция) — регламент и требования к отчёту полностью аналогичны очной защите;
- или асинхронную (почта, домашние задания ОРИОКС) — отчёт необходим, но оценивается защита; «**прислать лабы на почту**» без защиты = «неуд».

«На одном занятии» дистанционным группам читать как «за одну неделю из 1–16», «на следующем занятии» — как «через две недели».

Студенты очных групп, согласно требованиям МИЭТ, должны учиться очно. В виде исключения преподаватель лабораторных работ может принять у кого-то защиту дистанционно, но не обязан это делать.

### РЛ.6. Защита двух, трёх или более лабораторных работ (1–16 недели)

Две лабораторные работы на одном занятии могут быть оценены отдельно (как описано в разделе РЛ.4); если есть очередь на защиту — после защиты первой команда отправляется в конец очереди.

#### Комплект (три и более работы сразу) — оценивается выборочно, без бонусов

Если команда приносит на занятие три и более лабораторных работы сразу — защита происходит выборочно: преподаватель проверяет несколько произвольно выбранных (из всех выполненных, в том числе бонусных) заданий разных работ, задаёт вопросы и даёт задания по произвольным темам.

В каждую графу «Л<sub>і</sub>» выставляются одинаковые баллы:

$$\begin{cases} 0 \leq s \leq 3, & \text{если в комплекте 3–4 работы (до 12 недели),} \\ 0 \leq s \leq 1, & \text{если работ} \geq 5 \text{ или неделя} \geq 13. \end{cases} \quad (\text{Л.1})$$

Бонусные баллы, которые рассчитываются для одной лабораторной работы, при защите комплекта не рассчитываются и не добавляются:  $s^{\text{bonus}} = 0$ .

#### Третья лабораторная работа без предупреждения — ждёт две недели

Если команда уже получила баллы за две работы (или за комплект из 3–4) на занятии — защита следующей не ранее следующего занятия или через две недели; регламент — по фактической дате защиты, а не по первой заявке.

### РЛ.7. Лабораторные работы на 17–18 неделях и в сессию

Раздельная защита лабораторных работ (с вычислением баллов, бонусов и штрафов к каждой работе отдельно) на 17–18 неделях и в сессию **невозможна**.

Если, вопреки требованиям регламента, команда принесёт на 17 неделе и позже любое количество лабораторных работ (от одной и более), преподаватель может:

- либо оценить их как один комплект: выборочно, на  $0 \leq s \leq 1$  и  $s^{\text{bonus}} = 0$ ;
- либо сразу послать такую команду писать Кр1 и Кр2.

Исключение — если команда сдавала лабораторные работы весь семестр, начиная с первого занятия, по графику и досдаёт тому же преподавателю *одну* последнюю работу — она оценивается как сделанная на 16 неделе.

# Лабораторная работа 0

## Представление данных в ЭВМ

Только для тех, кто не изучал ОЭВМ и Асм! Изучавшие — начинают ОТИК с Л1.

Максимальная оценка ПИН — 0 баллов, НБ — 16 баллов.

Успешная защита Л0 = допуск к Л2–Л6. К Л1 допуск не требуется.

Для должников ОЭВМ и Асм: Л0 + Кр0 + базовый вопрос = закрытие долга.

Для индивидуальщиков обсуждается индивидуально.

### Л0.1. Задание на лабораторную работу

Все задания сформулированы для C/C++, как для ЯВУ ОТИК по умолчанию. Если для работы с файлами ОТИК планируется использовать другой ЯВУ — в Л0.№1 обязательно использовать не C/C++, а средства и/или документацию выбранного ЯВУ для оценки диапазона  $N$ -битного значения в файле.

**Штраф — 2 балла за каждую** лабораторную работу, где чтение/запись в файл выполняется средствами не того ЯВУ, который исследован в Л0.№1.

В Л0.№2–Л0.№6 использовать выбранный ЯВУ *рекомендуется*.

#### Рекомендации

Во всех заданиях Л0 рекомендуется вывод с помощью *printf()*:

- а) библиотека *libc*, в том числе *printf()*, доступна и в GNU Assembler;
- б) формат вывода *printf()* компактен и влияет только на одно значение.

Не рекомендуется вывод в потоки (доступен только в C++), так как:

- а) вывод в заданном формате (а не по умолчанию) объёмен и запутан;
- б) большинство манипуляторов не прекращают своё действие до отмены, а в комбинации с другими дают неочевидные эффекты;
- в) придётся выполнять лишние преобразования, так как по умолчанию в поток и *char* / *unsigned char*, и *int8\_t* / *uint8\_t*, и указатели на них выводятся как символы/строки.

Тем не менее, вывод в потоки не штрафует. Часть заданий различается для *printf()* и для потоков из-за принципиально разного подхода к выводу.

**Задание Л0.№1.** При помощи оператора *sizeof* выясните, сколько байтов занимают переменные типов: *char*, *unsigned char*, *short*, *unsigned short*, *int*, *unsigned*, *long long*, *unsigned long long*, *float*, *double*.

Для каждого типа рассчитайте разрядность в битах, учитывая, что для архитектуры x86/amd64 байт = октет = 8 бит.



**Штраф – 2 балла**, если выводятся только числа, без пояснений, и непонятно, где размер какого типа. **Бонус + 2 балла**, если при помощи макроса пояснения выводятся так, что в коде каждое имя типа в **Л0.№1** встречается единожды.

Напечатайте, используя `<climits>` (`limits.h`) С или `<limits>` C++, минимальные и максимальные значения исследуемых знаковых и беззнаковых целых типов.

1. Сколько различных значений может принимать переменная беззнакового  $N$ -битного типа? Знакового  $N$ -битного?  
Связано ли это как-то со значением  $N$  и как именно?
2. Каждое ли целое число  $x \in [0, 2^N)$  имеет своё  $N$ -битное представление? Всякая ли последовательность  $\xi$  из  $N$  битов может быть рассмотрена как целое  $x \in [0, 2^N)$ ? Взаимно однозначно ли это соответствие?
3. Каждое ли целое число  $x \in [-2^{N-1}, 2^{N-1})$  имеет своё  $N$ -битное представление? Всякая ли последовательность  $\xi$  из  $N$  битов может быть рассмотрена как целое  $x \in [-2^{N-1}, 2^{N-1})$ ? Взаимно однозначно ли это соответствие?

Напечатайте минимальные и максимальные значения `min` и `max` исследуемых типов с плавающей запятой.

1. Каждое ли вещественное число  $x \in [\min, \max]$  имеет своё представление с плавающей запятой стандарта IEEE 754 (*float/double* ЭВМ)?
2. Всякая ли последовательность  $\xi$  из  $N$  битов ( $N \in \{32, 64\}$ ) может быть рассмотрена как  $N$ -битное значение с плавающей запятой?  
Всегда ли это значение — число?
3. Каких чисел больше: 32-битных целых (*int/unsigned*) или 32-битных с плавающей запятой (*float*)? 64-битных целых (*long long/unsigned long long*) или 64-битных с плавающей запятой (*double*)?

Значения  $\pm\infty$  числами не являются, нечисла `nan` — тем более.

Если используется отличный от C/C++ язык — исследуйте:

- знаковые и беззнаковые целые числа разрядности 8, 16, 32 и 64;
- числа с плавающей запятой одинарной и двойной точности.

Если в языке нет средств для работы с числами заданной разрядности, а не только с заабстрагированными до полной переносимости «целочисленными вообще» и «с плавающей запятой вообще» — язык не подходит для ОТИК, меняйте.

Если есть и то, и другое — чтение и запись числа в файл всегда должны иметь конкретные, известные автору и обоснованные разрядность и формат.

**Задание Л0.№2.** Исследуйте внутреннее представление 16-битных чисел.

Для этого на C/C++ разработайте `void print16(void *p)` (одну: С или П).

С-**Л0.№2** для вывода при помощи стандартной библиотеки С (рекомендуется): функция `print16()` интерпретирует  $p$  как адрес 16-битного целого числа  $x$  типа *short/unsigned short* и печатает  $x$  при помощи `printf()`:

С-а) в шестнадцатеричном представлении;

С-б) в двоичном представлении: если не поддерживается формат  $b$  — напечатать биты  $(x \& (1 \ll i)) \neq 0$ , от старшего  $i = 15$  и до  $i = 0$ ;

С-в) в десятичном беззнаковом представлении;

С-г) в десятичном знаковом представлении.

Необходимый вид вывода `print16()` для значений 13 и 0x8000 по адресу  $p$ :

```
000D 0000000000001101    13    +13
```

```
8000 1000000000000000 32768 -32768
```

Для `print16()`, а также последующих `print32()`, `print64()`, и любых  $x$  и  $y$  младшая цифра любого представления  $y$  всегда должна печататься под младшей цифрой соответствующего представления  $x$ .

**П-10.№2** для вывода в потоки (крайне не рекомендуется, но и не штрафуются).

Так как вывод в поток различает *short* и *unsigned short*, функция `print16()` интерпретирует адрес  $p$  дважды:

- как адрес 16-битного беззнакового целого  $ux$  типа *unsigned short*;
- и как адрес 16-битного знакового целого  $sx$  типа *short*;

и печатает оба  $ux$  и  $sx$  во всех доступных представлениях:

П-а)  $ux$  в шестнадцатеричном представлении;

П-б)  $ux$  в двоичном представлении (шаблон `std::bitset<N>`);

П-в)  $ux$  в десятичном представлении;

П-г)  $sx$  в шестнадцатеричном представлении;

П-д)  $sx$  в двоичном представлении (шаблон `std::bitset<N>`);

П-е)  $sx$  в десятичном представлении.

Не забывайте, что манипуляторы *dec/hex* действуют на все числа до отмены, а `setw(int w)` — только на одно следующее.

Убедитесь в процессе исследования, что (П-б) и (П-д) — одно и то же двоичное представление; (П-а) и (П-г) — одно и то же шестнадцатеричное представление. Если у вас (П-б)  $\neq$  (П-д) или (П-а)  $\neq$  (П-г) — ищите ошибку.

Сократите вывод П-10.№2 до вида, аналогичного С-10.№2.

**Штраф – 1 балл**, если версия задания С-10.№2 реализована при помощи потоков или П-10.№2 при помощи стандартной библиотеки С.

**Штраф – 1 балл**, если вывод `print16()` занимает более одной строки.

Исследуйте при помощи `print16()` 16-битные целочисленные переменные типа *short/unsigned short*, принимающие значения:

- минимальное и максимальное целое беззнаковое 16-битные значения;
- минимальное и максимальное целое знаковое 16-битные значения;
- целочисленные  $x$ ,  $y$ ,  $a$ ,  $b$ , соответствующие варианту (таблица 10.1).

Как представляются в памяти ЭВМ знаковые целые значения? Как различаются целочисленные  $x$  и  $y = -x$ ?

Варианты значений

Таблица Л0.1

$(N^0 - 1)\%2 + 1$	Вариант
1	$x = 9, y = -9, a = 1, b = 2, c = 12345678, d = 123456789$
2	$x = 5, y = -5, a = 1, b = 2, c = 12345689, d = 123456891$

Как связаны двоичное представление (С-б) и шестнадцатеричное (С-а)? Как по шестнадцатеричному (С-а) записать двоичное для любого выбранного числа?

**Задание Л0.№3.** Исследуйте внутреннее представление 32-битных чисел — целых *int/unsigned* и с плавающей запятой *float*. Для этого на С/С++ разработайте `void print32(void *p)` (также одну— С или П — соответственно Л0.№2).

- С-Л0.№3 для стандартной библиотеки С: `print32()` интерпретирует *p* дважды:
- как адрес 32-битного целого числа *x* типа *int/unsigned*;
  - как адрес 32-битного числа с плавающей запятой *fx* типа *float*;
- и печатает *x* в представлениях (С-а)–(С-г), а *fx*:
- С-д) в шестнадцатеричном экспоненциальном представлении;
- С-е) в десятичном экспоненциальном представлении;
- С-ж) в представлении с десятичной запятой.

```
Необходимый вид вывода print32() для значений 13 и 0x80000000 по адресу p:
0000000D 00000000000000000000000000000000 13 +13
+0X1.A0P-146 +1.82e-44 +0.00
80000000 10000000000000000000000000000000 2147483648 -2147483648
-0X0.00P+0 -0.00e+00 -0.00
```

Если ширина терминала позволяет вывести (С-а)–(С-ж) в одну строку — это необходимо сделать; иначе вторая строка должна иметь отступ.

- П-Л0.№3 для вывода в потоки: `print32()` интерпретирует адрес *p* трижды:
- как адрес 32-битного беззнакового целого *ux* типа *unsigned*;
  - как адрес 32-битного знакового целого *sx* типа *int*;
  - как адрес 32-битного числа с плавающей запятой *fx* типа *float*;
- и печатает *ux* и *sx* в представлениях, аналогичных (С-а)–(С-г) (дублирующиеся — в одном экземпляре), а *fx* — в представлениях, аналогичных (С-д)–(С-ж). Обратите внимание, что манипулятор `setprecision(int n)` действует до отмены; `fixed/scientific/hexfloat` — также до отмены, а также проявляют себя по-разному в зависимости от `dec/hex`.

Исследуйте при помощи `print32()` 32-битные переменные — целочисленные типа *int/unsigned* и с плавающей запятой типа *float*:



Чем различается одно и то же значение ( $x$  или  $a$ ), записанное в *float* и *double*? В каком интервале *double*-чисел больше:  $[+0, +2)$  или  $[+2, +\infty)$ ?

**Штраф — 4 балла**, если для любой *printWW()*, где  $WW \in \{16, 32, 64\}$ :

- порядок представлений отличается от указанного;
  - представления не соответствуют друг другу (а команда, вместо того, чтобы задать вопрос преподавателю — так и сдаёт);
  - в частности, двоичное (C-6) явно несообразно шестнадцатеричному (C-a);
- или формат вывода отличается от указанного в худшую сторону, в частности:
- количество выводимых цифр двоичного представления (C-6) отлично от количества бит  $WW$  в числе;
  - а выводимых цифр шестнадцатеричного (C-a) — от количества тетрад  $\frac{WW}{4}$ ;
  - младшая цифра не под младшей предыдущего вызова *printWW()*;
  - вывод менее компактен, чем указано (делать компактнее при сохранении информативности — можно и нужно).

Далее во всех заданиях **всех лабораторных работ**, если явно не указано иное, все числа (как целые, так и с плавающей запятой) разрядности  $WW$  — как **результат**, так и **исходные данные** — должны выводиться соответствующей *printWW*  $\in \{\text{print16}(), \text{print32}(), \text{print64}()\}$ .

**Штраф — 2 балла за каждое число**, выведенное только в одном представлении (— 4 балла за число, если это единственное представление — десятичное).

**Задание Л0.№5.** Разработайте функцию *c16to32(void \*p)*, которая принимает адрес 16-битной целочисленной переменной, печатает её значение *print16()*, расширяет это значение до 32 бит двумя способами:

- как знаковое (*short*  $\rightarrow$  *int*);
- как беззнаковое (*unsigned short*  $\rightarrow$  *unsigned int*).

и для каждого способа печатает результат *print32()*.

Явное расширение в C++ выполняется *static\_cast*; неявное — в частности, при присваивании. Если источник и приёмник различаются не только размером, но и знаковостью — неопределённое поведение C/C++.

Проверьте её работу на значениях  $m$  и  $n$  (таблица Л0.2).

**Задание Л0.№6.** Разработайте функцию *ab16(void \*p)*, которая принимает адрес 16-битной целочисленной переменной  $x$  и выполняет над её копиями операции (a1)–(b6). Оригинал, лежащий по адресу  $p$ , должен оставаться неизменным; для каждой операции исходным являється значение  $x$ , а не результат предыдущей. Исходное значение  $x$  и каждый результат печатается *print16()*.

Сопоставьте результаты (a*i*) и (b*i*) — вначале на значении  $x = m$ , затем  $x = n$  (таблица Л0.2). Сколько всего различных операций описано в задании Л0.№6?

Варианты целочисленных значений

Таблица Л0.2

$(N^0 - 1) \% 3 + 1$	Вариант
1	$m = 57, n = -21$
2	$m = 21, n = -37$
3	$m = 37, n = -33$

- а1) беззнаковое умножение на 2;

а2) знаковое умножение на 2;

а3) беззнаковое деление на 2;

а4) знаковое деление на 2;

а5) расчёт остатка от беззнакового деления на 16;

а6) округление вниз до числа, кратного 16 (беззнаковое);
- б1) беззнаковый сдвиг влево на 1 бит;

б2) знаковый сдвиг влево на 1 бит;

б3) беззнаковый сдвиг вправо на 1 бит;

б4) знаковый сдвиг вправо на 1 бит;

б5) расчёт  $x \ \& \ 15$ ;

б6) расчёт  $x \ \& \ -16$ .

Если Л0.№6 реализуется на чистом C/C++, то беззнаковые и знаковые операции (умножение/деление/сдвиги) записываются одним и тем же оператором. Таким образом, адрес  $p$  необходимо интерпретировать дважды:

- как адрес 16-битного беззнакового целого  $ux$  типа *unsigned short*;
- и как адрес 16-битного знакового целого  $sx$  типа *short*.

Одни и те же знаки операций C/C++:

- для  $ux$  обозначают беззнаковые операции (если операция бинарная, типа  $*$  или  $/$ , то и второй операнд должен быть того же типа *unsigned short*);
- для  $sx$  — знаковые (бинарные — для  $sx$  и другого *short*).

Также можно забежать вперёд и реализовать операции как *ассемблерные вставки* — на уровне ассемблера беззнаковые и знаковые операции выполняются разными командами: *mul/imul, shl/sal, div/ldiv, shr/sar*.

Л0.2. Дополнительные бонусные и штрафные баллы

- 2 балла за каждое задание, где не печатаются исходные данные.
- 3 балла за утечку памяти (выделенные, но не освобождённые блоки).

# Лабораторная работа 1

## Просмотр и редактирование файлов в шестнадцатеричном представлении. Работа с файлами в C/C++

Максимальная оценка ПИН — 7 баллов.

Штраф за одно пропущенное обязательное задание — 2 балла.

### Л1.1. Задание на лабораторную работу

**Задание Л1.№1.** С помощью hexdump/xxd или шестнадцатеричного просмотрщика/редактора исследуйте файлы различных форматов (некоторые файлы представлены в папке «labs-files/Файлы в разных форматах»). Выделите сигнатуры или иные признаки формата там, где это возможно.

Описания части форматов находятся в папке «labs-files/Описание некоторых форматов», для прочих можно найти в Сети.

**Задание Л1.№2.** Определите тип файла, соответствующего номеру варианта ((№ – 1)%10), из папки «labs-files/Варианты 1 — \*». Откройте его корректным приложением. Поместите в отчёт тип и описание содержимого файла, а также признаки формата, по которым удалось определить тип.

**Задание Л1.№3.** В соответствии с номером варианта отредактируйте (таблица Л1.1) изображение colorchess16x16x2.bmp, используя шестнадцатеричный редактор. Откройте изменённый файл и убедитесь, что изменения корректны.

Файл colorchess16x16x2.bmp представляет собой изображение  $16 \times 16$  пикселей, сиренево-болотное (две сиреневые и две болотные клетки по  $8 \times 8$  пикселей), глубина цвета — 1 бит на пиксель. Убедитесь, что просмотрщик отображает его как цветное (некоторые игнорируют палитру файлов с глубиной 1 бит на пиксель).

### Варианты действий для редактирования

Таблица Л1.1

$(\text{№} - 1) \% 3 + 1$	Вариант
1	Поставить зелёную точку в правом нижнем углу изображения
2	Поставить сиреневую точку в правом верхнем углу изображения
3	Поставить зелёную точку в левом верхнем углу изображения

**Задание Л1.№4.** С помощью hexdump/xxd или шестнадцатеричного просмотрщика/редактора исследуйте файлы формата «простой текст» (plain text), представленные в различных кодировках (папка labs-files/Файлы в формате простого текста - кодировки разные, раздел Л1.2).

Есть ли у простого текста заголовок?

Сравните один и тот же текст, представленный в различных кодировках: размер и шестнадцатеричное представление.

Сравните шестнадцатеричное представление с тем, как текстовый редактор читает этот текст в кодировке платформы (обычно UTF-8). Учтите, что малофункциональные редакторы, такие как Блокнот MS Windows, поддерживают только одну-две кодировки, и «кракозябры» для остальных — нормальное явление.

Обратите внимание на файлы Алфавит - \*, включающие 192 символа национальных кодовых таблиц русского языка / кодовой таблицы Unicode, в том числе:

- 189 печатных (пробелы, цифры, латинские и русские буквы, символ @);
- 3 управляющих (переводы строки LF в стиле UNIX).

Одинаково ли количество байтов для представления одного символа кодовой таблицы (печатного или управляющего, как LF) в различных кодировках?

Одинаково ли шестнадцатеричное представление пробела в различных кодировках? Цифр? Латинских букв? Русских букв?

## Л1.2. Кодировки и кодовые таблицы русского языка

### Кодировки и кодовые таблицы

Строго говоря, необходимо различать понятия:

- *коддовая таблица* или таблица кодов — соответствие символов кодам в каком-то диапазоне;
- *кодировка* — представление кода символа в памяти или на диске.

Но обычно их смешивают и называют для краткости «кодировкой» совокупность кодовой таблицы и собственно кодировки.

Это в большинстве случаев не приводит к недопониманию, так как:

- кодировкам UTF-8, UTF-16, UTF-32 всегда соответствует одна и та же универсальная кодовая таблица Unicode, содержащая *все* национальные алфавиты;
- национальным кодовым таблицам (KOI8-R, IBM CP866, Windows-1251 и т. п.) всегда соответствует одна и та же кодировка: код записывается октетом «как есть».

Таким образом, *однобайтовыми кодировками* на практике называются *коддовые таблицы*, сопоставляющие символы кодам в диапазоне 0–255 (0x00–0xFF); коды символов таких таблиц всегда записывается одним октетом (байтом x86). В настоящее время все такие кодовые таблицы сопоставляют кодам 0–127 те же символы, что и кодовая таблица ASCII (являются расширениями кодовой таблицы ASCII). Коды в диапазоне 128–255 описывают национальные кодовые страницы.

*Коддовая таблица ASCII* сопоставляет символы кодам в диапазоне 0–127 (0x00–0x7F) (см. приложение В). При этом понятие «однобайтовая кодировка ASCII» *не определено*, и в зависимости от контекста может описывать как исторические способы записи семибитных ASCII-кодов восьмью битами (старший бит



мог использоваться для контроля чётности, дублировать один из семи младших или всегда быть нулевым), так и Latin-1 (ISO 8859-1), и, чаще, ту однобайтовую кодировку, которая используется на компьютере говорящего.

*Кодовая таблица Unicode* сопоставляет символы кодам 0–0x10FFFF (кодам 0–127 соответствуют символы ASCII) то есть не может быть представлена однобайтовой кодировкой. Не всем Unicode-кодам соответствуют символы; так, коды 0xD800–0xDFFF зарезервированы для представления суррогатных пар в UTF-16.

Для кодирования символов Unicode используются три основные кодировки.

1. UTF-8, кодировка переменной длины (изначально от 1 до 6 октетов, позже ограничили до 4) для узких строк:

- 0xxx xxxx — ASCII-символы (коды от 0 до 127 = 0x7F) представляются одним байтом, равным ASCII-коду (старший бит — 0).

Прочие, включая 128–255, представляются не менее чем двумя байтами. Старшие биты первого из  $k$  байтов содержат  $k$  единиц, затем следует разделитель 0, затем — старшие биты Unicode-кода. Старшие биты последующих байтов — 10, так что представление символа в UTF-8 не равно его коду Unicode:

- 110x xxxx 10xx xxxx — символы с Unicode-кодами от 128 = 0x80 до 0x7FF, в том числе греческие, русские и арабские буквы, представляются двумя байтами;
- 1110 xxxx 10xx xxxx 10xx xxxx — далее символы до 0xFFFF, в том числе основные китайские и японские иероглифы — три;
- 1111 0xxx 10xx xxxx 10xx xxxx 10xx xxxx — прочие символы (с запасом — до 0x1F FFFF) могут быть представлены четырьмя байтами.

UTF-8 — наиболее распространённая в настоящее время кодировка Unicode.

2. UTF-16, кодировка переменной длины (от 1 до 2 элементов) для широких строк из двухоктетных (16-битных) элементов:

- символы с Unicode-кодами 0x0000–0xFFFF записываются одним 16-битным элементом «как есть»; символов с кодами 0xD800–0xDFFF не существует;
- символы 0x1 0000–0x10 FFFF — двумя элементами: первый лежит в диапазоне 0xD800–0xDBFF, второй — 0xDC00–0xDFFF (суррогатной парой);

UTF-16 — старейшая кодировка Unicode (первые версии Unicode уместались в один 16-битный элемент). UTF-16, в отличие от UTF-8 и UTF-32, не позволяет записать коды свыше 0x10 FFFF.

3. UTF-32, кодировка постоянной длины (один 32-битный элемент с запасом вмещает все Unicode-коды) для широких строк из четырёхоктетных (32-битных) элементов.

Кодировки UTF-16 и UTF-32 имеют варианты, соответствующие разному порядку байтов в двух- или четырёхоктетном элементе (LE/BE); по умолчанию подразумевается порядок байтов платформы (LE для x86).

В MS Windows «Unicode» обозначает устаревшую версию кодировки UTF-16 (включающую только 16-битные символы, но не суррогатные пары), что неверно. Однобайтовая кодировка в MS Windows (для русского языка используется кодовая страница Windows-1251) обозначается «ANSI».

### Представление русского языка

В папке labs-files/Файлы в формате простого текста - кодировки разные представлены основные кодировки/таблицы для русского языка:

1. Многобайтовые кодировки универсальной кодовой таблицы Unicode:
  - UTF-8 — суффикс имени файла utf8;
  - UTF-16 — суффикс utf16;
  - UTF-32 — суффикс utf32.
2. Однобайтовые расширения ASCII:
  - КОИ-8 (код обмена информацией 8-битный) для русского алфавита (KOI8-R), использовавшаяся в России до широкого распространения MS DOS/MS Windows — суффикс koi8r;
  - ISO 8859-5, разработанная ISO и IEC и одно время считавшаяся стандартной, но не использовавшаяся — суффикс iso;
  - IBM CP866 (альтернативная кодировка ГОСТ), использовавшаяся в русифицированной MS DOS — суффикс dos;
  - Windows-1251 (CP1251), используемая в русифицированной MS Windows — суффикс windows;
  - MacCyrillic, используемая в Mac OS X — суффикс maccyrillic.

Из-за совпадения кодов наиболее частотных русских букв утилиты распознавания кодировок регулярно путают Windows-1251 и MacCyrillic; для различения этих кодировок необходим дополнительный анализ.

### Л1.3. Вопросы

1. Для чего нужен шестнадцатеричный редактор?
2. Какие функции libc используются для чтения/записи бинарных файлов?
3. Известно, что файл содержит осмысленный русскоязычный текст в одной из представленных в данной работе кодировок. Можно ли, используя только шестнадцатеричный редактор, без частотного анализа, отличить UTF-8, UTF-16, UTF-32: а) друг от друга, б) от однобайтовой кодировки? По каким признакам?

## Лабораторная работа 2

### Исследование статистических характеристик исходных текстов (как бинарных файлов, так и файлов в формате простого текста). Работа с кодовыми таблицами русского языка

Максимальная оценка ПИН обязательной части — 7 баллов. Максимум графы 7 + 7 баллов включает бонусное задание **Л2.№4**.

Штраф за одно пропущенное обязательное задание — 3 балла.

#### Л2.1. Задание на лабораторную работу

**Задание Л2.№1.** Разработайте программу или используйте набор программ (в частности, это может быть набор скриптов-однострочников, использующий стандартные утилиты GNU/Linux), который по заданному файлу  $Q$  рассчитывает:

- длину  $n$  файла  $Q$  в символах первичного алфавита  $A_1$ ;
  - $\text{count}(a_j)$  — общее количество вхождений каждого из символов  $a_j \in A_1$  в  $Q$  (ненормированную целочисленную частоту  $\nu_{\text{ненорм}}(a_j)$ );
- и оценивает, строя модель источника символов по файлу  $Q = c_1 \dots c_n$  в виде *источника без памяти* (модель для сжатия без учёта контекста):
- вероятность  $p_{\text{БП}}(a_j) = \frac{\text{count}(a_j)}{n}$  каждого из символов  $a_j \in A_1$ ;
  - количество информации  $I_{\text{БП}}(a_j) = -\log_2 \left( p_{\text{БП}}(a_j) \right)$  [бит] в каждом  $a_j \in A_1$ ;
  - суммарное количество информации  $I_{\text{БП}}(Q) = \sum_{i=1}^n I_{\text{БП}}(c_i)$  [бит],  $c_i \in Q$  в файле  $Q$  (не среднее на символ, которое в  $n$  раз меньше, а именно суммарное!).

Символом, как всегда, является *байт* (для x86/amd64 это 8 бит, *октет*), первичным алфавитом  $A_1 = \{a_1, \dots, a_{|A_1|}\}$  — множество возможных значений байта  $\{00, \dots, FF\}$  или, что то же самое,  $\{0, \dots, 255\}$ .

Сравните:

- а) длину файла  $Q$  в битах  $n \cdot 8$  и оценку  $I_{\text{БП}}(Q)$  в битах (печатайте оценку  $I_{\text{БП}}(Q)$  [бит] с двумя знаками после запятой, а также её дробную часть  $\{I_{\text{БП}}(Q) \text{ [бит]}\}$  в экспоненциальной форме);
- б) длину файла  $Q$  в октетах (8-битных байтах x86/amd64)  $n$  и оценку в октетах:
  - $I_{\text{БП}}(Q) \text{ [октетов]} = \frac{I_{\text{БП}}(Q) \text{ [бит]}}{8}$  (с двумя знаками после запятой);а также *оценки снизу* длин в октетах для сжатия без учёта контекста:
  - длины  $E = \lceil I_{\text{БП}}(Q) \text{ [октетов]} \rceil$  только сжатого текста, без информации, необходимой для его декодирования;

- длины  $G_{64} = E + 256 \cdot 8$  архива, где к сжатому тексту добавляется таблица  $(\nu(00), \dots, \nu(\text{FF}))$  из 256 ненормированных 64-битных частот  $\nu_{\text{ненорм}}(j) = \text{count}(j)$ ;
- длины  $G_8 = E + 256 \cdot 1$  архива, где к сжатому тексту добавляется таблица из 256 нормированных 8-битных частот  $\nu_{\text{норм}}(j)$ .

Все полученные величины необходимо напечатать на экране или сохранить в виде текстового файла-отчёта; причём таблица характеристик символов алфавита  $j \in A_1$  (символ, то есть байт  $j$  — в виде шестнадцатеричного значения байта, а не ASCII-символа; ненормированная частота  $\text{count}(j)$ ; оценки вероятности  $p(j)$  и количества информации  $I(j)$ ) должна либо печататься дважды:

- отсортированной по алфавиту (по значению  $j$ );
- отсортированной по убыванию  $\text{count}(j)$ ;

либо в программе необходимо предусмотреть пересортировку.

Проверьте разработанную программу на файлах различного формата (не только простом тексте; в том числе и на бинарных).

Выгодно ли применять к проанализированным файлам сжатие без учёта контекста? Нужно ли нормировать частоты?

Для побайтового чтения используйте  $fread()$  или её аналог в используемом языке программирования.

Для проверки корректности используйте файлы с заранее известным  $I_{\text{БП}}(Q)$ : для четырёх разных произвольных октетов  $a, b, c, d$  верно  $I_{\text{БП}}(a) = 0$  бит,  $I_{\text{БП}}(ab) = 2$  бита,  $I_{\text{БП}}(abcd) = 4 \cdot 2 = 8$  бит (1 байт x86) и т. п.

**Задание Л2.№2. Бонус +3 балла для пар и НБ, обязательное для ПИН-троек.** Разработайте программу, аналогичную Л2.№1, но считающую символом кодирования *печатный или управляющий символ* *Unicode*, а первичным алфавитом  $A_1$  — множество символов Unicode (строчных букв, заглавных букв, цифр, различных пробельных символов, знаков препинания и т. п.) в файле  $Q$ .

Обратите внимание, что Л2.№2 и связанное с ним Л2.№5 — *единственные* задания курса ОТИК, где символ кодирования не обязательно является байтом, а под исходным текстом не всегда понимается произвольный бинарный файл.

При оценивании длины архива учитывайте, что первичный алфавит  $A_1$  Л2.№2, в отличие от Л2.№1, имеет переменные длину и состав. Соответственно, сохранять в архиве нужно не массив из 256 частот  $(\nu(00), \dots, \nu(\text{FF}))$ , а:

- 64-битную длину  $|A_1|$  алфавита;
- массив из  $|A_1|$  пар символ-частота:  $((a_1, \nu(a_1)), \dots, (a_{|A_1|}, \nu(a_{|A_1|})))$ , где либо в каждой паре символ  $a_j \in A_1 \subseteq \text{Unicode}$  имеет переменную длину (UTF-8 или UTF-16) либо в каждой паре 32-битен (UTF-32).

Какой будет длина в октетах этих данных для исследуемых файлов?

Какой была бы длина в октетах, если бы сохранялись частоты не для символов файла  $A_1 \subseteq \text{Unicode}$ , а для всего Unicode (длина и состав алфавита постоянны, сохраняются только частоты  $(\nu(0), \dots, \nu(|\text{Unicode}| - 1))$ )?

Исследуйте один и тот же длинный текстовый файл в кодировке UTF-8 программами Л2.№1 и Л2.№2. Как выбор первичного алфавита влияет на:

- оценку  $I_{\text{БП}}(Q)$  без учёта контекста;
- оценку снизу длин архива для сжатия без учёта контекста?

Повторите анализ для других длинных текстовых файлов в кодировке UTF-8.

**Задание Л2.№3.** Рассчитайте, используя программу Л2.№1, частоты октетов в файлах, являющихся простым текстом в различных кодировках (папка labs-files/Файлы в формате простого текста - кодировки разные). Исследуйте несколько разных осмысленных русскоязычных текстов и все представленные кодировки. Определите 4 наиболее частых октета среди всех используемых и 4 наиболее частых октета, не являющихся кодами печатных символов ASCII. Обратите внимание на распределение октетов многобайтовых кодировок.

Рассчитайте частоты октетов в файле, соответствующем варианту  $N \bmod 9$  в папке labs-files/Варианты 2 - определение кодировки простого текста (далее — файл  $W$ ).

Определите, является ли  $W$  простым русскоязычным текстом в одной из стандартных кодировок (один из вариантов представляет собой нерусскоязычный текст); если да — определите кодировку.

**Задание Л2.№4. Бонус +7 (добавляется к л/р, а не к общему бонусу) — источник с памятью.**

Разработайте программу, которая по заданному файлу  $Q$  рассчитывает:

- $\text{count}(a_j a_k)$  — количество вхождений подстрок  $a_j a_k$ ;
- $\text{count}(a_j *)$  — общее количество вхождений любых двухсимвольных подстрок, начинающихся с  $a_j$  (для всех символов, кроме последнего символа файла,  $\text{count}(a_j *) = \text{count}(a_j)$ );

и оценивает, строя модель источника символов по файлу  $Q = c_1 \dots c_n$  в виде *стационарного источника Маркова первого порядка*:

- условную вероятность  $p(a_k | a_j) = \frac{\text{count}(a_j a_k)}{\sum_k \text{count}(a_j a_k)} = \frac{\text{count}(a_j a_k)}{\text{count}(a_j *)}$  каждой пары символов  $a_j, a_k \in A_1$ ;
- суммарное количество информации  $I_{\text{СМ1}}(Q)$  в файле  $Q$  в битах и байтах (безусловную вероятность  $p(c_1 = a_j)$  считайте равной  $\frac{1}{256}$ );

аналогично Л2.№1 (символ кодирования — байт, как всегда).

Проверьте разработанную программу на файлах различного формата (не только простом тексте; в том числе и на бинарных). Сопоставьте результат с Л2.№1.

Таблицу частот какого размера нужно сохранить вместе со сжатым текстом длины  $I_{\text{СМ1}}(Q)$ , чтобы успешно его декодировать? Выгодно ли такое сжатие?

**Задание Л2.№5. Бонус +3 балла.** Разработайте программу, аналогичную Л2.№4, но считающую символом кодирования *печатный или управляющий символ Unicode*, аналогично Л2.№2.

Проверьте разработанную программу, сопоставьте результат с Л2.№4 и Л2.№2.

## Л2.2. О терминах

**Символ кодирования** в ТИ есть элемент (квант) качественной информации. В кодировании (то есть сжатии, защите от помех или шифровании) символом является, как правило, **байт** (для x86 — октет); изредка — отдельный бит (символ вторичного алфавита Хаффмана или первичного алфавита Хэмминга) или битовый блок отличной от байта длины; и никогда — печатный символ ASCII, KOI-8, Unicode etc (технически реализовать можно, как показывают задания Л2.№2/Л2.№5, но смысла в этом нет).

**Текст, строка** — последовательность символов в указанном выше смысле (в общем случае бинарный файл и его фрагменты). **Алфавит** — множество всех возможных символов в указанном выше смысле.

Не обманывайтесь принятой в кодировании терминологией!

Кто обманется и воспользуется функцией *fgetc()* (или, тем более, *readLine()*, *split()* и т. п.) где-то, кроме Л2.№2/Л2.№5 — минус 2 балла.

## Лабораторная работа 3

### Проектирование формата файлов

**Цель работы:** научиться создавать и обрабатывать двоичные файлы на ЯВУ.

Максимальная оценка ПИН обязательной части — 7 баллов. Максимум графы  $7 + 7$  включает бонус за архивацию папок.

Штраф за одно пропущенное обязательное задание — 2 балла.

#### ЛЗ.1. Задание на лабораторную работу

**Задание ЛЗ.№1.** Выберите сигнатуру для собственного формата файлов, которая будет использоваться во всех дальнейших работах (6 байт).

Разработайте программу-кодек, состоящую из двух частей — *кодера* и *декодера*.

1. *Кодер* по заданному файлу  $Q$  создаёт архив  $R$  формата, описанного в таблице ЛЗ.1.

#### Нулевая версия формата архива

Таблица ЛЗ.1

Смещение поля (байты)	Размер поля (байты)	Описание поля
0	6	Сигнатура формата (выбранное выше 6-байтовое значение)
6	2	Версия формата (беззнаковое 16-битное целое): для задания ЛЗ.№1 — 0
8	8	Исходная длина $n$ файла в байтах (беззнаковое 64-битное целое)
16	$n$	«Сырые» данные исходного файла (несжатый текст)

2. *Декодер* по заданному архиву  $R$ :

- проверяет сигнатуру на соответствие выбранной в ЛЗ.№1 и версию формата;
- при корректных сигнатуре и версии восстанавливает файл  $\tilde{Q}$  (который должен совпадать с исходным  $Q$ ) по данным архива  $R$ .

Кодер и декодер могут быть реализованы как в виде двух отдельных модулей, так в одной программе (в последнем случае действие задаётся ключом командной строки или выбором меню — то есть должна быть возможность отдельного вызова кодера и декодера, а не только слитного преобразования  $Q \rightarrow R \rightarrow \tilde{Q}$ ).

Проверьте при помощи шестнадцатеричного просмотрщика/редактора, что поля заголовка созданного архива соответствуют таблице Л3.1; при помощи утилиты GNU/Linux *diff* проверьте побайтовое совпадение распакованного  $\bar{Q}$  с исходным  $Q$ .

При выполнении работы в MS Windows аналог *diff* можно найти, используя справочную команду консоли *help*.

**Задание Л3.№2.** Разработайте и опишите формат файла-архива для дальнейших работ (не обязан развивать таблицу Л3.1; формат Л3.№2 может отличаться от Л3.№1 размером сигнатуры, но её начало должно сохраниться). Архив обязательно содержит заголовок и закодированные данные.

Заголовок обязательно включает:

- сигнатуру (в дальнейшем все дальнейшие версии архивов команды должны использовать сигнатуру из Л3.№2);
  - ненулевую версию формата (1 и выше; при разделении на мажорную и минорную — 0.1 и выше);
  - коды использованных алгоритмов сжатия и защиты от помех (с учётом того, что сжатие без учёта контекста применяется поверх сжатия с его учётом — кодов алгоритма сжатия будет два);
  - исходную длину файла в символах кодирования (то есть байтах);
- а также любые необходимые, по мнению автора, поля.

Формат должен предусматривать возможность добавления служебных данных для различных алгоритмов сжатия и защиты от помех (например, массив частот для методов сжатия без учёта контекста) без кардинальной его переработки.

В качестве кодов алгоритмов, соответствующих отсутствию сжатия и защиты от помех, используйте значение 0.

**Задание Л3.№3.** Разработайте программу-кодэк (аналогично Л3.№3), создающую архив формата, разработанного в задании Л3.№2.

## Л3.2. О терминах

**Символ кодирования = байт x86 = октет**, а не печатный символ ASCII, KOI-8, Unicode etc.

**Несжатый текст  $Q$  = любой бинарный файл**; сжатый текст, который следует в  $R$  после кода алгоритма и служебной информации — тоже бинарный файл.

Не обманывайтесь принятой в кодировании терминологией! (Кто обманется и воспользуется функциями *fgetc()* и т. п. — минус 2 балла).

«Маркером» текстового представления является не только *fgetc()*, но и — тем более — функции форматированного (что значит текстового) ввода-вывода и обработки строк, в частности, *readLine()*, *split()* и т. п.



В большинстве библиотек для ввода-вывода «сырых» бинарных данных есть только две функции:

- аналог *fread()* для чтения из файла;
  - и аналог *fwrite()* для записи;
- и их *достаточно*.

### ЛЗ.3. Дополнительные бонусные и штрафные баллы

−2 балла за текстовое представление архива или если кодер обрабатывает только файлы в текстовом представлении — см. раздел ЛЗ.2.

−2 балла, если чтение/запись в файл выполняется средствами не того ЯВУ, который исследован в ЛО.№1.

−1 балл, если структура архива, создаваемого кодеком в задании ЛЗ.№3, не совпадает с описанной в задании ЛЗ.№2.

−1 балл, если декодер при некорректной сигнатуре архива  $R$  всё равно создаёт файл  $\tilde{Q}$ .

+7 баллов, если в структуре заголовка ЛЗ.№2 и программе ЛЗ.№3 предусмотрена возможность собрать в один архив и восстановить ещё и иерархическую структуру папок (не более +4, если можно собрать несколько файлов, но пути не сохраняются).

## Лабораторная работа 4

### Сжатие данных без учёта контекста (энтропийное сжатие)

Максимальная оценка ПИН обязательной части — 14 баллов. Максимум графы 14 + 14 включает Л4.№5. Для НБ — бонус +16 полная, +8 упрощённая.

Штраф за одно пропущенное обязательное задание — 3 балла.

#### Л4.1. Упрощённое задание (не более 7 баллов за работу)

Используйте демонстрационные программы любой свободной библиотеки для сжатия методом Хаффмана, Шеннона—Фано, Шеннона либо целочисленным арифметическим (интервальным) методом. Узнайте из документации, какой метод используется и что является символом кодирования.

Продемонстрируйте работу кодера и декодера. Сопоставьте длину сжатых данных с оценкой количества информации в исходном файле (Л2.№1).

Отличие упрощённого задания от полного, реализованного с помощью библиотеки со свободной лицензией — в управлении заголовком.

В упрощённом задании студент использует демонстрационную программу «как есть»; полное может быть засчитано только в том случае, если запись/чтение заголовка реализованы студентом и могут быть изменены.

Засчитывается только в том случае, если студент может доказать, используя документацию библиотеки, что используется именно заявленный студентом метод сжатия без учёта контекста и только он (а не, например, связка LZW+Хаффман).

#### Л4.2. Задание на лабораторную работу

**Задание Л4.№1.** Разработайте программу-кодировщик (аналогично Л3.№3), реализующую сжатие/разжатие файла методом Хаффмана.

Сигнатура формата должна совпадать с выбранной в Л3.№1. Если в формат необходимо внести какие-либо изменения по сравнению с Л3.№2 — они должны быть описаны в отчёте, а номер версии — изменён.

Также в отчёте обязательно должны быть описаны:

- код алгоритма, не равный 0 (если выполняются бонусные задания ниже — разным алгоритмам должны соответствовать разные коды);
- состав информации для декодирования и формат её хранения в файле;
- принятые при построении дерева уточнения (порядок сортировки при совпадении частот и т. п.).

Сопоставьте длину сжатых данных с оценкой количества информации в исходном файле по стационарной модели без памяти (Л2.№1), а длину архива — с её оценкой снизу (Л2.№1).

**Задание Л4.№2.** Разработайте программу, рассчитывающую для файла коды Хаффмана (собственно архивы со сжатыми данными можно не создавать):

- 64) по ненормализованным частотам байтов  $\nu_{64}(j) = \text{count}(j)$  в диапазоне  $0 \dots 2^{64} - 1$  (для хранения одного значения необходимо  $|\nu_{64}| = 8$  байт);
- 32) по усечённым (для файла длиной до  $2^{32} - 1$ ) или нормализованным к диапазону  $0 \dots 2^{32} - 1$  для больших файлов частотам ( $|\nu_{32}| = 4$  байта);
- 8) по частотам, приведённым к диапазону  $0 \dots 255$  ( $|\nu_8| = 1$  байт);
- 4) по частотам, приведённым к диапазону  $0 \dots 15$  ( $|\nu_4| = \frac{1}{2}$  байта).

Частоты нормализовывать таким образом, чтобы ноль переходил в 0, единица в 1 (то есть чтобы ненулевые частоты не превращались в нулевые).

В отличие от Л2.№1, где рассчитывается оценка снизу длины сжатого текста и длины архива — здесь рассчитывается точное значение.

Для каждого  $|\nu_B|$  программа должна рассчитывать:

- длину  $E_B$  сжатых данных файла в байтах;
- общую длину  $G_B$  данных, необходимых для распаковки (сжатых данных  $E_B$  и массива частот) в байтах:  $G_B = E_B + 256 \cdot \frac{B}{8} = E_B + 32 \cdot B$ .

Для файла программа должна рассчитывать наиболее выгодную разрядность  $B^*$  частот  $\left(G_{B^*} = \min_B(G_B)\right)$  для сжатия методом Хаффмана.

Для каждого из нескольких различных файлов сравните  $E_{64}$ ,  $E_{32}$ ,  $E_8$  и  $E_4$ , а также  $G_{64}$ ,  $G_{32}$ ,  $G_8$  и  $G_4$ ; сравните  $B^*$ ; запишите в отчёт.

Насколько выводы Л4.№2, учитывающие округление частот при приведении к заданному диапазону, отличаются от предварительных выводов Л2.№1?

Выгодно ли встраивать в реализацию алгоритма подбор  $B^*$  для конкретного файла? Учитывайте время и необходимость сохранения такого  $B^*$  в заголовке.

Какая фиксированная для алгоритма разрядность  $B^{**}$  лучше подходит для использования (удобнее в чтении/записи и при этом даёт достаточно малое  $G_B$ )? Предложите способ представления  $256 B^{**}$ -битных частот в виде  $32 \cdot B^{**}$  байтов.

+2 балла, если программа перебирает для файла все возможные разрядности  $B$  от 1 до 64 бит и рассчитывает  $B^*$  среди всех, а не только 64/32/8/4.

+2 балла за графическое изображение зависимостей  $\frac{E_B}{E_{64}}$  и  $\frac{G_B}{G_{64}}$  от  $B$  для набора из 10 или более файлов разного типа. В итоге должно получиться одно изображение для  $E$  (либо со множеством графиков на нём, либо с «ящичками с усами») и одно для  $G$ , независимо от количества файлов — если на каждый файл будет свой график, баллы не начисляются.

**Задание Л4.№3.** Разработайте универсальный декодер разработанного формата, который:

- проверяет сигнатуру на соответствие выбранной в Л3.№1, при некорректной сигнатуре прекращает работу;
- при корректной сигнатуре — проверяет номер версии формата: для предыдущих версий формата вызывает свою старую версию, при ещё неизвестном номере версии прекращает работу;
- при корректных сигнатуре и версии — коды алгоритмов, после чего вызывает соответствующий им декодер из заданий Л3.№3 или Л4.№1.

В дальнейшем необходимо будет дополнять анализ при добавлении нового алгоритма (при этом любая модификация — например, реализация метода Хаффмана с другим порядком сортировки при совпадении частот — должна оформляться именно как новый алгоритм с новым кодом).

**Задание Л4.№4. Бонус +3 балла.**

Разработайте «интеллектуальный» кодер, который анализирует суммарный объём  $n_{comp}$  сжатых данных и информации для декодирования Л4.№1 и, если  $n_{comp} \geq n$ , записывает в полученный архив несжатый текст исходного файла (то есть использует вместо кодера Л4.№1 кодер Л3.№3; код алгоритма при этом также должен быть 0).

Сразу предусмотрите флаг для отключения анализа и принудительного использования заданного алгоритма.

**Задание Л4.№5. Бонус +14 (добавляется к л/р, а не к общему бонусу) — арифметический/интервальный кодек.**

Разработайте программу-кодек (аналогично Л4.№1–Л4.№4), реализующую сжатие/разжатие файла целочисленным арифметическим (интервальным) методом.

**Задание Л4.№6. Бонус +3 балла для пар, обязательное для троек.**

Разработайте программу-кодек (аналогично Л4.№1–Л4.№4), реализующую сжатие/разжатие файла методом, соответствующим варианту (таблица Л4.1).

**Варианты Л4.№6**

Таблица Л4.1

$(N^a - 1) \% 2 + 1$	Вариант
1	метод Шеннона
2	метод Шеннона—Фано

### Л4.3. Дополнительные бонусные и штрафные баллы

—4 балла за текстовое представление архива или если кодер обрабатывает только файлы в текстовом представлении — см. раздел Л3.2.

—2 балла, если массив частот может для какого-то одного файла иметь длину более 256 байтов.

—2 балла, если формат архива не соответствует Л3.№2 (штраф не начисляется, если в отчёте описана новая версия формата, которая логичнее исходной).

—4 балла, если даже сигнатура не совпадает с Л3.№1.

—2 балла, если ранее разработанный декодер Л3.№3 пытается декодировать архив Л4.№1, а не выдаёт сообщение о несоответствии алгоритма.

+1 балл, если при работе с несколькими файлами/папками каждый из файлов имеет собственный код алгоритма и информацию для декодирования (а в Л4.№4 — анализируется отдельно).

+4 балла, если архив позволяет как включить для каждого файла собственный код алгоритма и информацию для декодирования, так и рассмотреть совокупность файлов как единый исходный текст (но при декодировании восстановить исходную структуру файлов).

## Лабораторная работа 5

### Сжатие данных с учётом контекста

Максимальная оценка ПИН обязательной части — 14 баллов. Максимум графы 14 + 14 включает **Л5.№3**. Для НБ — бонус +16 полная, +8 упрощённая.

Штраф за одно пропущенное обязательное задание — 3 балла.

#### Л5.1. Упрощённое задание (не более 7 баллов за работу)

Разработайте кодер и декодер «наивного» RLE. Продемонстрируйте работу кодера и декодера.

**Штраф — 3 балла** за текстовое представление архива или если кодер обрабатывает только файлы в текстовом представлении — см. раздел **Л3.2**.

Засчитывается только в том случае, если все студенты команды знают:

- а) сколько байтов они отводят на код  $\tilde{L}$  длины цепочки  $L$ ;
- б) в каком виде записывают длину цепочки  $L \geq 1$  в эти байты:  $\tilde{L} = L - 1$  или  $\tilde{L} = L$ ;
- в) какое максимально допустимое значение  $L_{\max}$  для длины цепочки  $L$  следует из сделанного ими выбора (а) и (б);
- г) в каком порядке записывают код цепочки «наивного» RLE:  $(\tilde{L}, c)$  или  $(c, \tilde{L})$ ; и если они корректно обрабатывают большее, чем  $L_{\max}$ , количество повторений подряд одного символа  $c$  в файле.

#### Л5.2. Задание на лабораторную работу

**Задание Л5.№1.** Разработайте программу-кодек, реализующую сжатие/разжатие файла методом RLE согласно варианту (таблица **Л5.1**). Здесь и далее в разных вариантах используются разные коды/алгоритмы, реализующие один заданный метод — семейство кодов и алгоритмов. Если вы не поняли, какой код нужно реализовать в вашем варианте — спрашивайте. «Какая-нибудь» реализация заданного метода, не соответствующая варианту — не засчитывается.

Вычислите минимальные и максимальные возможные значения  $L$  для своего варианта (для всех возможных случаев).

Сопоставьте длину сжатых данных с оценкой количества информации в исходном файле по модели Маркова (**Л2.№4**), а также с объёмом таблицы частот в **Л2.№4**.

**Задание Л5.№2.** Разработайте программу-кодек, реализующую сжатие/разжатие файла методом LZ78 согласно варианту (таблица **Л5.3**).

**Задание Л5.№3. Бонус +14 (добавляется к л/р, а не к общему бонусу).** Разработайте программу-кодек, реализующую сжатие/разжатие файла методом LZ77 согласно варианту (таблица **Л5.3**).

## Варианты кодов семейства RLE

Таблица Л15.1

$(N^0 - 1)\%3 + 1$	Вариант
1	флаг-бит сжатая/несжатая цепочка; цепочка из $L \geq 3$ одинаковых символов $c...c$ как $(1:1, (L-3):7, c:8)$ , из $L \geq 1$ разных $c_1...c_L$ — как $(0:1, (L-1):7, c_1:8, \dots c_L:8)$
2	префикс $p$ ; цепочка из $L \geq 4$ одинаковых символов $c...c$ , $c \neq p$ как $(p:8, (L-3):8, c:8)$ , цепочка из $L \geq 2$ символов $p...p$ как $(p:8, (L-1):8, p:8)$ , одиночный $p$ как $(p:8, 0:8)$
3	префикс $p$ ; цепочка из $L \geq 4$ одинаковых символов $c...c$ , $c \neq p$ как $(p:8, (L-4):8, c:8)$ , цепочка из $L \geq 1$ символов $p...p$ как $(p:8, (L-1):8, p:8)$ , в том числе одиночный $p$ как $(p:8, 0:8, p:8)$

## Варианты кодов семейства LZ78

Таблица Л15.2

$(N^0 - 1)\%2 + 1$	Вариант
1	концепт Зива и Лемпеля 1978 года
2	LZW

Запись в таблице  $(S:10, L:6)$  обозначает: сначала записывается 10 бит  $S$ , затем 6 бит  $L$ . Для концепта Зива и Лемпеля символ  $c$  записывается как  $(c:8)$  (то есть просто как байт  $c$ ); для реализации с односимвольным префиксом  $p$  символ  $c \neq p$  также записывается как  $(c:8)$ . Если ссылка отделяется от несжатого текста односимвольным префиксом  $p$ , значение  $p$  выбирать для каждого исходного текста  $Q$  отдельно, исходя из частот символов в  $Q$ .

Вычислите минимальные и максимальные возможные значения  $S$  и  $L$  для своего варианта.

**Задание Л15.№4. Бонус +3 балла для пар, обязательное для троек.** До-работайте универсальный декодер из задания Л14.№3 и кодер из задания Л14.№4 с учётом всех реализованных алгоритмов.

Варианты кодов семейства LZ77

Таблица Л5.3

$(N^0 - 1) \% 7 + 1$	Вариант
1	концепт Зива и Лемпеля; ссылка $\{S, L\}$ как $(S : 10, L : 6)$
2	концепт Зива и Лемпеля; ссылка $\{S, L\}$ как $(L : 6, S : 10)$
3	флаг-бит ссылка/цепочка; ссылка $\{S, L\}$ как $(1 : 1, (L - 3) : 7, (S - 1) : 8)$ , цепочка $c_1 \dots c_L$ как $(0 : 1, (L - 1) : 7, c_1 : 8, \dots c_L : 8)$ ,
4	флаг-биты ссылка/символ группируются во флаг-байты; ссылка $\{S, L\}$ записывается как $(S : 10, (L - 3) : 6)$ — LZJB
5	флаг-биты ссылка/символ группируются во флаг-байты; ссылка $\{S, L\}$ записывается как $((L - 3) : 6, S : 10)$
6	префикс $p$ ; ссылка $\{S, L\}$ как $(p : 8, S : 10, (L - 4) : 6)$ , символ $c = p$ в тексте — как $(p : 8, 0 : 8, 0 : 8)$
7	префикс $p$ ; ссылка $\{S, L\}$ как $(p : 8, (L - 3) : 6, (S - 1) : 10)$ , символ $c = p$ в тексте — как $(p : 8, 0 : 8)$

**Задание Л5.№5. Бонус +7 для пар, +3 для троек, обязательное для четвёрок.** Доработайте универсальные кодер и декодер так, что к одному файлу было возможно применить последовательно:

- 1) вначале сжатие с учётом контекста;
- 2) затем сжатие без учёта контекста.

Баллы не начисляются, если кодер/декодер сводятся к последовательному запуску двух программ, то есть если на первом проходе формируется заголовок формата Л3.№2, и второй проход рассматривает этот заголовок как часть исходного текста.

Конкретные алгоритмы выбираются при запуске: сжатие с контекстом из реализованных в данной л/р, без контекста — из реализованных в предыдущей л/р.



## Лабораторная работа 6

### Помехозащитное кодирование

Исключена после переноса ОТИК с четвёртого года обучения на третий. Если кто-то выполнит — максимальная оценка обязательной части 14 баллов; оценка добавляется к сумме бонусов.

Штраф за одно пропущенное обязательное задание — 2 балла.

#### Л6.1. Упрощённое задание (не более 7 баллов за работу)

Используйте демонстрационные программы любой свободной библиотеки для защиты от помех методом Хэмминга либо Рида—Соломона. Узнайте из документации, какой метод используется, какое количество ошибок  $E$  он исправляет в блоке, размер блока, тип исправляемой ошибки.

Продемонстрируйте, что  $E$  ошибок при декодировании исправляется, а  $E + 1$  — уже нет.

#### Л6.2. Задание на лабораторную работу

**Задание Л6.№1.** Разработайте программу-кодек, реализующую защиту данных от помех методом Хэмминга, позволяющую исправить одиночную инверсию в блоке и обнаружить двойную инверсию в блоке.

Размер блока до кодирования соответствует варианту (таблица Л6.1). Код

#### Варианты размера блока до кодирования

Таблица Л6.1

$(N - 1) \% 2 + 1$	Вариант
1	$N = 7$ байт
2	$N = 8$ байт

должен быть систематическим, количество контрольных символов — целым (если в контрольном символе при заданном размере блока остаются неиспользуемые биты — продублируйте бит общей чётности).

Внесите в закодированные данные ошибки; убедитесь, что при декодировании действительно исправляется одиночная инверсия в блоке и обнаруживается двойная.

**Задание Л6.№2.** Разработайте программу-кодек, реализующую защиту данных от помех методом Рида—Соломона (используйте библиотеки, распространяемые по свободным лицензиям). Какова выбранная вами длина блока до кодирования (сколько информационных символов)? Сколько контрольных символов

добавляется? Какое максимальное количество ошибок в блоке может быть исправлено?

Внесите в закодированные данные ошибки; убедитесь, что при декодировании действительно исправляется указанное выше количество ошибок.

**Задание Л6.№3.** К ключевым полям заголовка файла (сигнатура, версия, алгоритмы и т. п.) нельзя применить какой-либо помехозащитный код без нарушения читаемости заголовка, поэтому они должны обрабатываться отдельно (например, троироваться, или дублироваться с применением к дубликату защиты от помех, или добавленные контрольные символы заголовка должны быть помещены вне заголовка).

Добавьте новую версию формата, где заголовок защищён от помех даже в том случае, когда к данным защита не применяется.

**Задание Л6.№4. Бонус +3 балла для пар, обязательное для троек.** Доработайте универсальный декодер из задания Л4.№3 и кодер из задания Л4.№4 с учётом реализованных алгоритмов.

**Задание Л6.№5. Бонус +7 для пар, +3 для троек, обязательное для четвёрок.** Доработайте универсальные кодер и декодер так, что к одному файлу было возможно применить последовательно:

- 1) вначале сжатие с учётом контекста;
- 2) затем сжатие без учёта контекста;
- 3) в конце защиту от помех.

Баллы не начисляются, если кодер/декодер сводятся к последовательному запуску двух программ, то есть если на первом проходе формируется заголовок формата Л3.№2, и второй проход рассматривает этот заголовок как часть исходного текста.

Дополнительно +7 баллов, если кодер и декодер включают шифрование:

- 1) вначале сжатие с учётом контекста;
- 2) затем сжатие без учёта контекста;
- 3) после всех алгоритмов сжатия — шифрование;
- 4) в конце — защиту от помех.

### Л6.3. Дополнительные бонусные и штрафные баллы

- 2 балла, если размер блока Хэмминга — весь файл.
- 2 балла, если код несистематический.
- 2 балла, если в блоке после кодирования есть неиспользуемый бит, а двойная ошибка в блоке не распознаётся.

Приложение В. Коды ASCII

ASCII CONTROL CODE CHART

b7 b6 b5 BITS  b4 b3 b2 b1				0		0		0		0		1		1		1		1									
				0		0		1		0		1		0		1		0		1							
				CONTROL						SYMBOLS NUMBERS						UPPER CASE						LOWER CASE					
0 0 0 0				0	NUL		16	DLE		32	SP		48	0		64	@		80	P		96	'		112	p	
0 0 0 1				1	SOH		17	DC1		33	!		49	1		65	A		81	Q		97	a		113	q	
0 0 1 0				2	STX		18	DC2		34	"		50	2		66	B		82	R		98	b		114	r	
0 0 1 1				3	ETX		19	DC3		35	#		51	3		67	C		83	S		99	c		115	s	
0 1 0 0				4	EOT		20	DC4		36	\$		52	4		68	D		84	T		100	d		116	t	
0 1 0 1				5	ENQ		21	NAK		37	%		53	5		69	E		85	U		101	e		117	u	
0 1 1 0				6	ACK		22	SYN		38	&		54	6		70	F		86	V		102	f		118	v	
0 1 1 1				7	BEL		23	ETB		39	'		55	7		71	G		87	W		103	g		119	w	
1 0 0 0				8	BS		24	CAN		40	(		56	8		72	H		88	X		104	h		120	x	
1 0 0 1				9	HT		25	EM		41	)		57	9		73	I		89	Y		105	i		121	y	
1 0 1 0				10	LF		26	SUB		42	*		58	:		74	J		90	Z		106	j		122	z	
1 0 1 1				11	VT		27	ESC		43	+		59	;		75	K		91	[		107	k		123	{	
1 1 0 0				12	FF		28	FS		44	,		60	<		76	L		92	\		108	l		124		
1 1 0 1				13	CR		29	GS		45	-		61	=		77	M		93	]		109	m		125	}	
1 1 1 0				14	SO		30	RS		46	.		62	>		78	N		94	^		110	n		126	~	
1 1 1 1				15	SI		31	US		47	/		63	?		79	O		95	`		111	o		127	DEL	
				F			17	1F		37	2F		57	3F		77	4F		117	5F		137	6F		157	7F	

LEGEND:

dec	CHAR	
hex		oct

Victor Eijkhout  
Dept. of Comp. Sci.  
University of Tennessee  
Knoxville TN 37996, USA