

COMPUTER SCIENCE TRIPOS Part II

Thursday 8 June 2006 1.30 to 4.30

PAPER 9

Answer *five* questions.

Submit the answers in five *separate* bundles, each with its own cover sheet. On each cover sheet, write the numbers of *all* attempted questions, and circle the number of the question attached.

You may not start to read the questions
printed on the subsequent pages of this
question paper until instructed that you
may do so by the Invigilator

STATIONERY REQUIREMENTS

Script Paper

Blue Coversheets

Tags

1 Human–Computer Interaction

The development manager of a website for online book-buying has asked you to carry out a heuristic evaluation of its usability. He has specifically proposed the three heuristics listed below.

1. “There should be between five and nine navigation options on each page.”
2. “There should be a good match between the navigation buttons and the users’ goals.”
3. “It should be easy for users to change their plans.”

(a) Has the manager misunderstood heuristic evaluation? Briefly justify your answer. [2 marks]

(b) Please comment on the above three heuristics suggested by the manager. For *each* of the proposed heuristics, your comments should include:

- (i) any theoretical justification for (or against) this heuristic;
- (ii) any additional evaluation steps that might be required in applying it; and
- (iii) the likely impact of such evaluation on the system design.

[6 marks each]

2 VLSI Design

(a) What is *dual-rail logic*? Why is it useful in self-timed circuits? [4 marks]

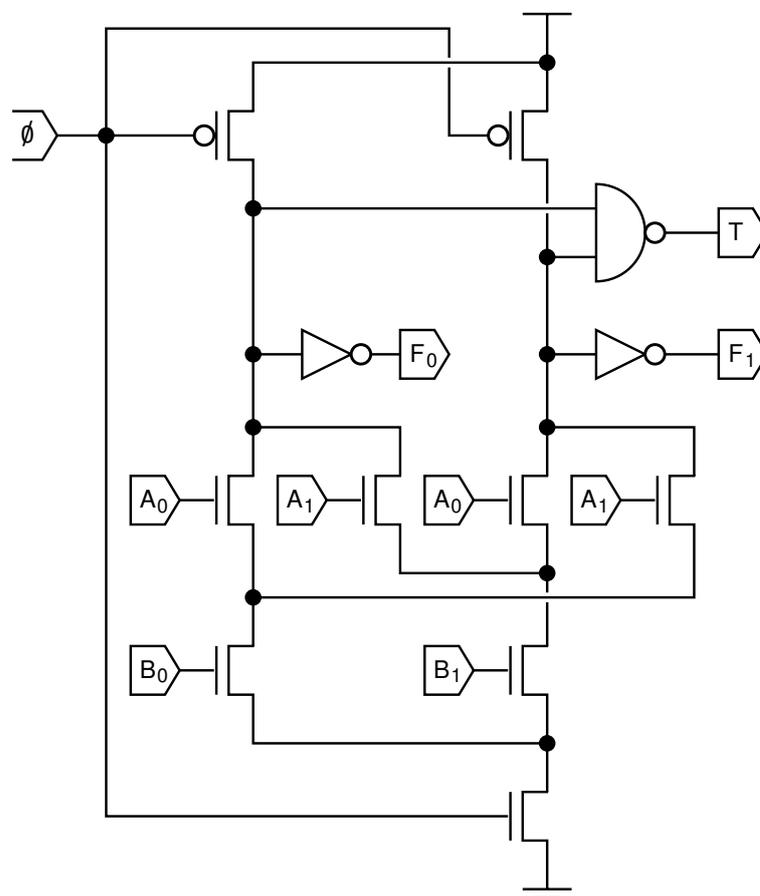
(b) Using dual-rail logic, sketch circuits for

(i) an inverter; [2 marks]

(ii) an AND gate; [3 marks]

(iii) an exclusive OR gate. [3 marks]

(c) The following circuit shows a dynamic dual-rail gate:



What function does it implement? Explain how it works. [8 marks]

3 Digital Communication II

- (a) Switches are used in a variety of different communications networks.
- (i) Describe how switch architectures that are used to forward samples in Time Division Multiplexed networks, such as the digital telephone network, are designed to scale to large numbers of input and output ports. [5 marks]
 - (ii) In asynchronous networks such as Asynchronous Transfer Mode (ATM) or Internet Protocol (IP) packet networks, switches are also used to forward packets. What are the basic components of a router with input and output buffering? [5 marks]
- (b) Distributed routing algorithms in communications systems are designed to provide a fault-tolerant computation of end-to-end paths in the event of link or router failure (or repair).
- (i) Describe how this occurs, using as an example the distance-vector algorithm. [5 marks]
 - (ii) Distance-vector routing is said to be slow to react to changes. Explain why, and outline why link-state protocols are therefore preferred in today's Internet. [5 marks]

4 Distributed Systems

Members of an open process group manage distributed replicas of data values stored in persistent memory. To allow the system to operate in the presence of transient failures of some replica managers, a quorum assembly scheme is used. Replica managers are assumed to be non-malicious and fail-stop.

To update a managed data item, the operations provided by the managing process include:

```
lock(item)
update(item, value, timestamp)
read(item, timestamp)
unlock(item)
```

- (a) Suppose the data item is an initially empty list of values and the update operation appends a value. Illustrate the quorum assembly scheme for five replicas, showing a number of update and read operations. [8 marks]
- (b) How is a total order of updates achieved by quorum assembly in the presence of concurrent update requests by clients to the open group? Discuss how any problems that might arise can be solved. [4 marks]
- (c) When can `unlock(item)` be executed safely by the initiating replica manager? Describe any additional protocol that is needed. [5 marks]
- (d) Suppose that the process group is managing non-overlapping partitions of a distributed database instead of replicas. Can quorum assembly play any part in making the related updates required for distributed transactions? Justify your answer. [3 marks]

5 Advanced Systems Topics

- (a) Distributed storage approaches can be divided into *network attached storage* (NAS) and *storage area networks* (SANs). Explain with the aid of a diagram the basic differences between the two approaches. [4 marks]
- (b) The network file system (NFS) is often used in local area networks.
- (i) Why is NFS not normally considered suitable for wide area networks? [2 marks]
- (ii) Briefly discuss how one could modify NFS to better support wide area networks. [2 marks]
- (c) *Distributed shared virtual memory* can be used within a computing cluster to transparently allow multi-threaded programs to run across multiple machines. Sketch the design of a DSVM system. Be sure to explain what happens both when a memory read and when a memory write occurs. Comment on the expected performance and robustness of your system. [6 marks]
- (d) EROS is a capability-based operating system.
- (i) What is a capability? [1 mark]
- (ii) Explain with the aid of a diagram how EROS uses traditional paging hardware to emulate capability hardware. [5 marks]

6 Computer Vision

- (a) Extraction of visual features from images often involves convolution with filters that are themselves constructed from combinations of differential operators. One example is the Laplacian $\nabla^2 \equiv \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}$ of a Gaussian $G_\sigma(x, y)$ having scale parameter σ , generating the filter $\nabla^2 G_\sigma(x, y)$ for convolution with the image $I(x, y)$. Explain in detail each of the following three operator sequences, where $*$ signifies two-dimensional convolution.

(i) $\nabla^2 [G_\sigma(x, y) * I(x, y)]$ [2 marks]

(ii) $G_\sigma(x, y) * \nabla^2 I(x, y)$ [2 marks]

(iii) $[\nabla^2 G_\sigma(x, y)] * I(x, y)$ [2 marks]

(iv) What are the differences amongst their effects on the image? [2 marks]

- (b) In human vision, the photoreceptors responsible for colour (cones) are numerous only near the fovea, mainly in the central ± 10 degrees. Likewise high spatial resolution is only found there. So then why does the visual world appear to us uniformly coloured? Why does it also seem to have uniform spatial resolution? What implications and design principles for computer vision might be drawn from these observations? [4 marks]

- (c) Outline a scheme for accomplishing transcription of handwriting (not cursive, that is, with letters already separated). Explain the core modules in your system, from low-level feature extraction to high-level classification of letters. At the highest level of the classifier, explain how the system could use Bayesian methods to incorporate expert knowledge such as a lexicon of actual words and knowledge about relative letter frequencies. [4 marks]

- (d) How can dynamic information about facial appearance and pose in video sequences (as opposed to mere still-frame image information), be used in a face recognition system? Which core difficult aspects of face recognition with still frames become more tractable with dynamic sequences? Are some aspects made more difficult? [4 marks]

7 Advanced Graphics

Describe, in detail, the radiosity method for calculating illumination. Ensure that your answer gives an overview of the algorithm, describes an implementable method of calculating form factors, and explains an efficient way of iterating to a solution.

[20 marks]

8 Optimising Compilers

Consider the ML-like language given by abstract syntax

$$e ::= x \mid n \mid \lambda x.e \mid e_1 e_2 \mid \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \mid \text{store } e \text{ in } r \mid \text{load } e \text{ from } r$$

where x ranges over variable names, n over integer constants, and r over global names for disjoint areas of memory known as *regions*. This language allows values to be stored inside regions: *store e in r* writes the value of e at some newly allocated memory location within region r and returns a pointer to this new location; the complementary operation *load e from r* reads the value which e points to (provided e is indeed a pointer into region r , otherwise the operation fails without accessing r).

Types have syntax

$$\tau ::= \text{int} \mid \tau \rightarrow \tau \mid * \tau \text{ in } r$$

where $*\tau \text{ in } r$ is the type of a pointer to a τ -typed value stored in region r . Note that there is no polymorphism and that *if-then-else* uses an integer (rather than boolean) condition.

- (a) Give an *effect system* (also known as an annotated type system) in which we can derive judgements of the form

$$\Gamma \vdash e : t, \varphi$$

where t is an extended form of τ and Γ is a set of assumptions of the form $x : t$. Effects φ are sets of region names representing the regions which e may need to access (i.e. write into or read from) during its execution.

[12 marks]

- (b) Give types and effects for the following expressions, commenting briefly on any problems your scheme encounters and how they may be resolved. (Assume that r and s are region names, x is a variable of type $*\text{int in } r$, and p is a variable of type $*\text{int in } s$.)

(i) *if load x from r then store 42 in s else p* [2 marks]

(ii) *λy . if load x from r then store y in s else p* [2 marks]

(iii) *if load x from r then λy . store y in s else λy . p* [4 marks]

9 Artificial Intelligence II

In this question we deal with a general two-class supervised learning problem. Instances are denoted by $x \in X$, the two classes by c_1 and c_2 , and $h : X \rightarrow \{c_1, c_2\}$ denotes a hypothesis. Labelled examples appear independently at random according to the distribution P on $X \times \{c_1, c_2\}$. The loss function $L(c_i, c_j)$ denotes the loss incurred by a classifier predicting c_i when the correct prediction is c_j .

- (a) Show that if our choice of hypothesis h is completely unrestricted and L is the 0–1 loss function then the *Bayes optimal classifier* minimising

$$E[L(h(x), c)]$$

where the expected value is taken according to the distribution P is given by

$$h(x) = \begin{cases} c_1 & \text{if } \Pr(c_1|x) > \frac{1}{2} \\ c_2 & \text{otherwise.} \end{cases}$$

[10 marks]

- (b) We now define a procedure for the generation of training sequences, denoted by s . Let \mathcal{H} be a set of possible hypotheses, let $p(h)$ be a prior on \mathcal{H} , let $p(x)$ be a distribution on X and let $\Pr(c|x, h)$ be a likelihood, denoting the probability of obtaining classification c given instance x and hypothesis $h \in \mathcal{H}$. A training set s is generated as follows. We obtain a single $h \in \mathcal{H}$ randomly according to $p(h)$. We then obtain m instances (x_1, \dots, x_m) independently at random according to $p(x)$. Finally, these are labelled according to the likelihood such that

$$p(s|h) = \prod_{i=1}^m \Pr(c_i|x_i, h)p(x_i).$$

We now wish to construct a hypothesis h' , not necessarily in \mathcal{H} , for the purposes of classifying future examples. The usual approach in a Bayesian context would be to construct the hypothesis

$$h'(x) = \begin{cases} c_1 & \text{if } \Pr(c_1|x, s) > \frac{1}{2} \\ c_2 & \text{otherwise.} \end{cases}$$

By modifying your answer to part (a) or otherwise, show that this remains an optimal procedure in the case of 0–1 loss.

[10 marks]

10 Digital Signal Processing

Consider a software routine that converts the sampling rate of digital audio data from 8 kHz to 48 kHz, without changing the represented sound. It reads an input sequence $\{x_i\}$ and produces an output sequence $\{y_i\}$. The routine first inserts five samples of value 0 between each consecutive pair of input samples. This results in a new intermediate sequence $\{x'_i\}$ with $x'_{6i} = x_i$ and $x'_{6i+k} = 0$ for all $k \in \{1, \dots, 5\}$. The sequence $\{x'_i\}$ is then low-pass filtered, resulting in $\{y_i\}$.

- (a) How can the process of taking discrete-time samples $\{x_i\}$ from a continuous waveform $x(t)$ be modelled through a function $\hat{x}(t)$ that represents the sampling result but can still be analysed using the continuous Fourier transform? [2 marks]
- (b) What effect does sampling with 8 kHz have on the Fourier spectrum of the signal? [2 marks]
- (c) How and under what condition can this sampling process be reversed? [2 marks]
- (d) Can $\hat{x}(t)$ also model another sampling process that results in the discrete sequence $\{x'_i\}$, and if so, what is its sampling frequency? [2 marks]
- (e) How does the continuous spectrum associated with $\{x'_i\}$ relate to that of $\{x_i\}$? [2 marks]
- (f) What purpose is served by the low-pass filter that the routine applies? In particular, what would happen to a 1 kHz sine tone input if this filter were not applied and $\{y_i\} = \{x'_i\}$ were output instead? What cut-off frequency must the filter have? [5 marks]
- (g) Provide a formula for calculating a 25-sample long causal finite impulse response $\{h_i\}$ of a low-pass filter suitable for this routine, based on the Hamming windowing function. [5 marks]

11 Types

Given any polymorphic lambda calculus (PLC) type τ and any function ρ mapping type variables α to values $n \in \{-1, 0, 1\}$, a value $\llbracket \tau \rrbracket \rho$ in $\{-1, 0, 1\}$ is defined by recursion on the structure of τ as follows.

$$\llbracket \alpha \rrbracket \rho = \rho(\alpha)$$

$$\llbracket \tau_1 \rightarrow \tau_2 \rrbracket \rho = \begin{cases} 1 & \text{if } \llbracket \tau_1 \rrbracket \rho \leq \llbracket \tau_2 \rrbracket \rho \\ \llbracket \tau_2 \rrbracket \rho & \text{otherwise} \end{cases}$$

$\llbracket \forall \alpha (\tau) \rrbracket \rho =$ the minimum of the values $\llbracket \tau \rrbracket (\rho[\alpha \mapsto n])$ for $n = -1, 0, 1$ (where $\rho[\alpha \mapsto n]$ is the function mapping α to n and every other α' to $\rho(\alpha')$).

If Γ is a non-empty PLC typing environment, let $\llbracket \Gamma \rrbracket \rho$ denote the minimum value of $\llbracket \tau \rrbracket \rho$ as τ ranges over the types in Γ ; in the case that Γ is empty, we define $\llbracket \Gamma \rrbracket \rho$ to be 1.

- (a) Prove that if $\Gamma \vdash M : \tau$ is a valid PLC typing judgement, then for any ρ , $\llbracket \Gamma \rrbracket \rho \leq \llbracket \tau \rrbracket \rho$.

You may assume without proof that if α is not free in τ then

$$\llbracket \tau \rrbracket (\rho[\alpha \mapsto n]) = \llbracket \tau \rrbracket \rho$$

and also that type substitutions $\tau[\tau'/\alpha]$ satisfy

$$\llbracket \tau[\tau'/\alpha] \rrbracket \rho = \llbracket \tau \rrbracket (\rho[\alpha \mapsto \llbracket \tau' \rrbracket \rho])$$

[Hint: show that the property $\Phi(\Gamma, M, \tau) =$ “for all ρ , $\llbracket \Gamma \rrbracket \rho \leq \llbracket \tau \rrbracket \rho$ ” is closed under the rules of the typing system.]

[16 marks]

- (b) Deduce that there is no closed PLC expression of type

$$\forall \alpha, \beta (((\alpha \rightarrow \beta) \rightarrow \alpha) \rightarrow \alpha)$$

[4 marks]

12 Numerical Analysis II

- (a) With reference to solution of the differential equation $y' = f(x, y)$, explain the conventional notation $x_n, y(x_n), y_n, f_n$. [3 marks]
- (b) Explain the terms *local error*, *global error*, and *order* of a method. [3 marks]
- (c) Without deriving any formulae, describe the general technique for deriving *multistep* formulae. [2 marks]
- (d) Milne's method uses the multistep formulae

$$y_{n+1} = y_{n-3} + \frac{4h}{3}(2f_n - f_{n-1} + 2f_{n-2})$$

$$y_{n+1} = y_{n-1} + \frac{h}{3}(\tilde{f}_{n+1} + 4f_n + f_{n-1})$$

which each have local error $O(h^5)$. What is the meaning of the term \tilde{f}_{n+1} ? Suggest a suitable starting procedure and explain how the Milne formulae are used. [6 marks]

- (e) Let $x_0 = 0.2$, $y(x_0) = 1.67$, $h = 0.2$ and

$$f(x, y) = 1 + \frac{(y - x)(x + 2)}{x + 1}.$$

Suppose the following values of f_n have been generated by the starting procedure: 4.6, 5.6, 7.2 for $n = 1, 2, 3$. Calculate the first required value of \tilde{f}_{n+1} to 2 significant digits. [3 marks]

- (f) Contrast Milne's method with your starting procedure, commenting particularly on *stability*, *efficiency* and *step size* considerations. [3 marks]

13 Bioinformatics

- (a) Hidden Markov models (HMM) are widely used in Bioinformatics.
- (i) In a HMM when would you use the Baum–Welch algorithm, and when the Viterbi algorithm, and why? Give biologically motivated examples. [8 marks]
- (ii) Any machine learning model (such as a HMM) for protein secondary structure determination or gene finding relies on discovering characteristic statistical properties of protein sequences. Name a property (and justify your answer) that helps to localise (and distinguish) transmembrane segments and coils in a protein sequence, or exon/intron boundaries in a genomic region. [2 marks]
- (b) Discuss the complexity of an algorithm to reconstruct a genetic network from microarray perturbation data. [7 marks]
- (c) What is the difference in terms of connectivity between a scale-free network and a random network? Give biological examples of scale-free networks. [3 marks]

14 Denotational Semantics

- (a) State carefully, without proof, the soundness and adequacy results for PCF. (You need not describe the syntax, operational and denotational semantics of PCF in any detail.) [3 marks]
- (b) Define the logical relation you would use in proving adequacy for PCF. State carefully without proof the “fundamental theorem” for the logical relation. [5 marks]
- (c) Define contextual equivalence for PCF. [2 marks]
- (d) Explain carefully the difficulties in obtaining a fully abstract denotational semantics for PCF. [7 marks]
- (e) Describe how PCF syntax and operational semantics can be extended to obtain full abstraction. [3 marks]

15 Specification and Verification II

A JK flip-flop has inputs J, K and an output Q, which is driven by a stored value, and behaviour specified by the following table.

J	K	Q	Q _{next}
0	0	0	0
0	0	1	1
0	1	X	0
1	0	X	1
1	1	0	1
1	1	1	0

Assume that the stored value is initially 0.

- (a) Describe the sequence of outputs on Q if the J and K inputs are always 1. [2 marks]
- (b) Define a predicate JK such that $JK(j, k, q)$ models the behaviour of a JK flip-flop. Describe and justify the logical type of JK. [6 marks]
- (c) Write down a formal specification of a device COUNT such that the output at time t is $t \bmod 4$. [2 marks]
- (d) Design an implementation of COUNT using JK flip-flops, describe how it works and draw a diagram of your design. [4 marks]
- (e) How you would prove that your design meets its specification? [6 marks]

16 Topics in Concurrency

(a) Define what it means for two states in a transition system to be bisimilar. [2 marks]

(b) Hennessy–Milner logic has assertions

$$A ::= \bigwedge_{i \in I} A_i \mid \neg A \mid \langle \alpha \rangle A ,$$

where I is a set, possibly empty, indexing a collection of assertions A_i , and α ranges over a set of actions Act . Define the semantics of the logic within a transition system with actions in Act . [4 marks]

(c) Show that if two states in the transition system are bisimilar, then they satisfy the same assertions of the logic. [6 marks]

(d) Show that if two states in the transition system satisfy exactly the same assertions of the logic, then they are bisimilar. [8 marks]

END OF PAPER