# Mathematical Model of Computer Viruses

## PhD. Ferenc Leitold,

**Veszprém University - Veszprog Ltd., Hungary**

fleitold@veszprog.hu

VESZPROG

# Table of contents

- **Models of computation**
- **Operating system**
- **Virus definition**
- **What can we do with this mathematical model ?**

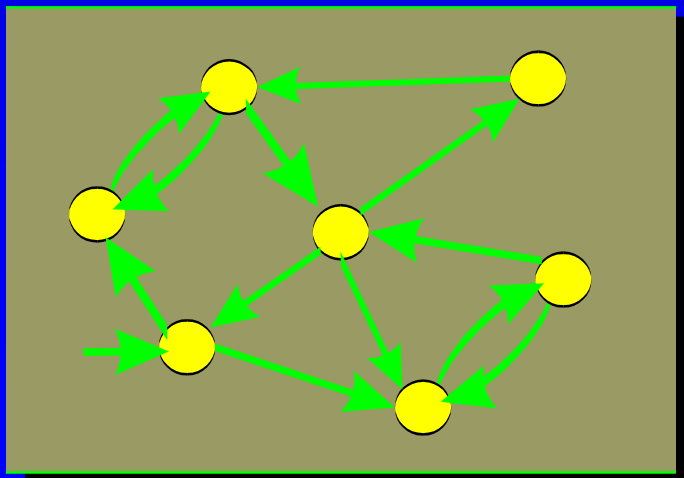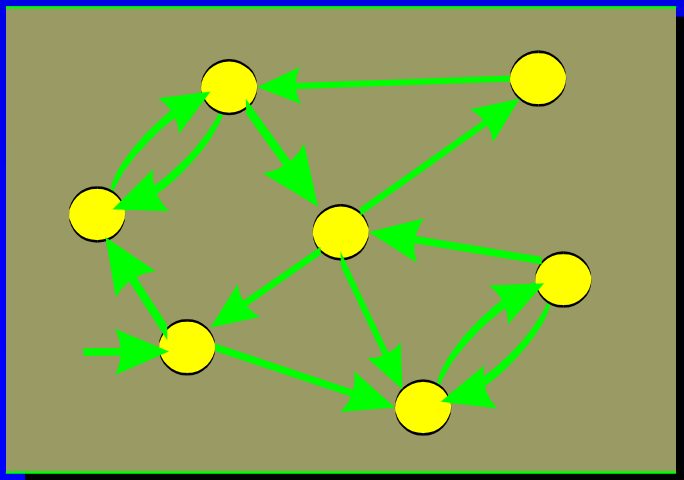# Turing Machine

# Turing Machine

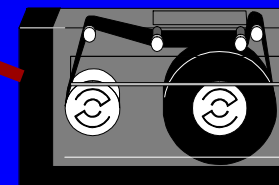**VESZPROG**

## Finite automata
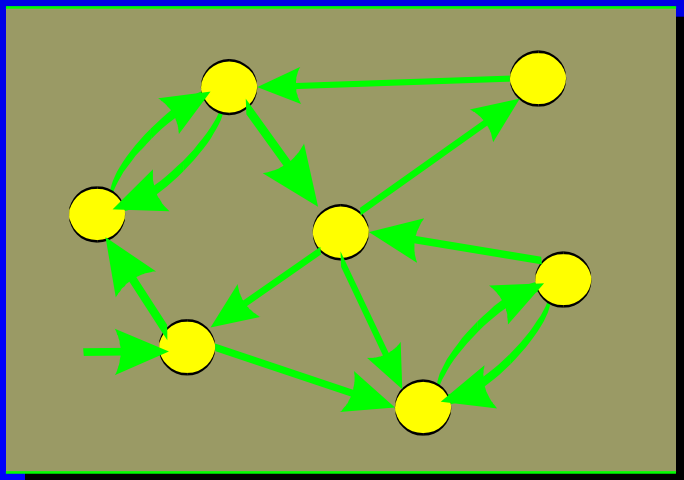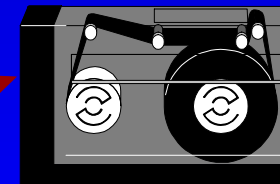
# Turing Machine

**Finite automata**

**Input tape**

# Turing Machine
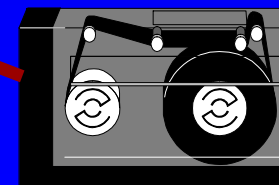
**Finite automata**

**Output tape**

**Input tape**

# Turing Machine

$$T = \langle Q, S, I, \delta, b, q_0, q_f \rangle$$

**S: tape symbols**
**I: input symbols, $I \subset S$**
**b: blank symbol, $b \in S \setminus I$**
**$\delta$: move function, $\delta : Q \times S \rightarrow Q \times S \times \{l, r, s\}$**

**Output tape**

**Finite automata**

Q

$q_f$

$q_0$

**Input tape**

Random Access Machine

Output tape

Memory

41256

VESZPROG

80486 DX4

Input tape

CPU & Program

# Random Access Machine

**Memory**

**CPU & Pro...**

| $m_0$ | **Accumulator** |
|-------|-----------------|
| $m_1$ | |
| $m_2$ | |
| $m_3$ | |
| $m_4$ | |
| $\vdots$ | |

# RASPM with SABS

Output tape

Memory
(Program)

Background tapes

Input tape

CPU

41256

80486
DX4

VESZPROG

# RASPM with ABS definition

**VESZPROG**

$$G = \langle V,U,T,f,q,M \rangle$$

M: initial memory content

q: initial value of the IP

$f : U \rightarrow T$

T: set of processor's activities

U: operation codes, $U \subseteq V$

V: set of symbols

# Instruction set

- **move** (LOAD, STORE)
- **logical** (AND, OR, XOR)
- **arithmetic** (ADD, SUB, MULT, DIV)
- **branch** (JUMP, JGTZ, JZERO)
- **input/output tape handling** (READ, WRITE)
- **background tape handling** (GET, PUT, SEEK, SETDRIVE)

VESZPROG

# Operating System

# Operating System

- **system of programs**

# Operating System

- **system of programs**
- **able to handle separate program or data files**

# Operating System

- system of programs
- able to handle separate program or data files
- able to make a specified program to run.

# Operating Systems
# under RASPM with ABS

VESZPROG

# Operating Systems under RASPM with ABS

- **The OS is in the initial memory (M)**

# Operating Systems under RASPM with ABS

- **The OS is in the initial memory (M)**
  - → **OS specific machine**

# Operating Systems under RASPM with ABS

- **The OS is in the initial memory (M)**
  - $\rightarrow$ **OS specific machine**
- **The OS is in the background tape**

# Operating Systems under RASPM with ABS

**VESZPROG**

- **The OS is in the initial memory (M)**
  → **OS specific machine**
- **The OS is in the background tape**
  → **OS independent machine**

# Operating Systems under RASPM with ABS

- **The OS is in the initial memory (M)**
  - **→ OS specific machine**
- **The OS is in the background tape**
  - **→ OS independent machine**
- **The OS is in the input tape**

# Operating Systems under RASPM with ABS

- **The OS is in the initial memory (M)**
  - → **OS specific machine**
- **The OS is in the background tape**
  - → **OS independent machine**
- **The OS is in the input tape**
  - → **unusable**

# Operating Systems under RASPM with ABS

- **The OS is in the initial memory (M)**
  → **OS specific machine**
- **The OS is in the background tape**
  → **OS independent machine**
- **The OS is in the input tape**
  → **unusable**

# Sample OS

# Comparing RASPM with ABS-es

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

# Comparing RASPM with ABS-es

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

$$\{q_1, M_1\} \neq \{q_2, M_2\}$$

VESZPROG

# Comparing RASPM with ABS-es

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

$$\{q_1, M_1\} \neq \{q_2, M_2\}$$

- **different operating systems**
- **different loader program**

# Comparing RASPM with ABS-es

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$
$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

**VESZPROG**

# Comparing RASPM with ABS-es

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

$$\{f_1, T_1, U_1\} \neq \{f_2, T_2, U_2\}$$

VESZPROG

# Comparing
# RASPM with ABS-es

**VESZPROG**

$$G_1 = \langle V_1, U_1, T_1, f_1, q_1, M_1 \rangle$$

$$G_2 = \langle V_2, U_2, T_2, f_2, q_2, M_2 \rangle$$

$$\{f_1, T_1, U_1\} \neq \{f_2, T_2, U_2\}$$

- **different activities**
- **different operation codes**

# Comparing RASPM with ABS-es

$$G_1 = \langle V_1, U_1, T_1, f_1, q_1, M_1 \rangle$$

$$G_2 = \langle V_2, U_2, T_2, f_2, q_2, M_2 \rangle$$

**VESZPROG**

# Comparing RASPM with ABS-es

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

$$V_1 \neq V_2$$

# Comparing RASPM with ABS-es

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

$$V_1 \neq V_2$$

- **different symbols**
- **different tape formats**

# Computer virus

# Computer virus

- a (part of) program

# Computer virus

- a (part of) program
- it is attached to a program area

# Computer virus

- a (part of) program
- it is attached to a program area
- it is able to link itself to other program areas

# Computer virus

- a (part of) program
- it is attached to a program area
- it is able to link itself to other program areas
- it is executed when the host program area is to be executed

# Virus spreading modes

# Virus spreading modes

- **machine specific**

# Virus spreading modes

- **machine specific**
- **machine independent**

# Virus spreading modes

- **machine specific**
- **machine independent**
- **operating system specific**

# Virus spreading modes

- **machine specific**
- **machine independent**
- **operating system specific**
- **operating system independent**

# Virus spreading modes

- **machine specific**
- **machine independent**
- **operating system specific**
- **operating system independent**

- **direct**

# Virus spreading modes

- **machine specific**
- **machine independent**
- **operating system specific**
- **operating system independent**

- **direct**
- **indirect**

# Virus spreading modes

- machine specific
- machine independent
- operating system specific
- operating system independent
- direct
- indirect

# Sample virus

# What can we do with this mathematical model ?

# What can we do with this mathematical model ?

- **Examine the working mechanism of viruses**

# What can we do with this mathematical model ?

**VESZPROG**

- Examine the working mechanism of viruses

- Examine the virus detection problem

# What can we do with this mathematical model ?

**VESZPROG**

- Examine the working mechanism of viruses

- Examine the virus detection problem

- Examine multiplatform viruses

# General virus detection problem

**VESZPROG**

**Theorem:**

**It is impossible to build a Turing Machine which could decide if an executable file in a RASPM with ABS contains a virus or not.**

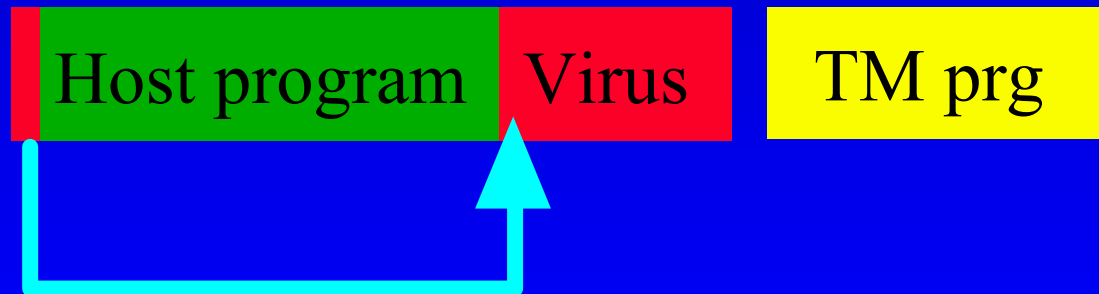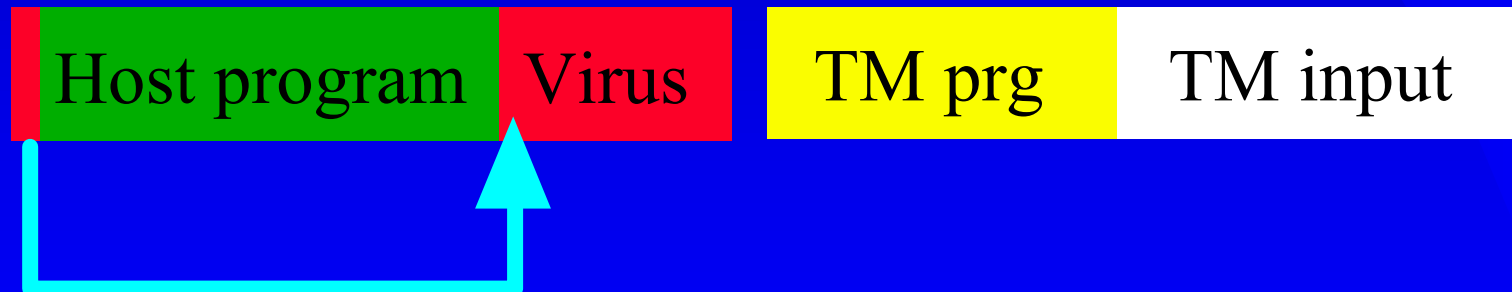# General virus detection problem

## Proof:

Host program | Virus

# General virus detection problem

**Proof:**

| Host program | Virus | TM prg |

# General virus detection problem

**VESZPROG**

## Proof:

| Host program | Virus | TM prg | TM input |

# General virus detection problem

## Proof:

| Host program | Virus | TM prg | TM input |

# General virus detection problem

## Proof:

# General virus detection problem

**Proof:**



| Host program | Virus | TM prg | TM input |

**Virus detection problem ➡ TM halting problem**

VESZPROG

# Multiplatform viruses

$$G_1 = \langle V_1, U_1, T_1, f_1, q_1, M_1 \rangle$$

$$G_2 = \langle V_2, U_2, T_2, f_2, q_2, M_2 \rangle$$

# Multiplatform viruses

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

Conditions:

$$V_1 \, ☋ \, U_2 \neq 0$$

$$U_1 \, ☋ \, V_2 \neq 0$$

$G_1$ has to know some operation codes of $G_2$

$G_2$ has to know some operation codes of $G_1$

VESZPROG

# Multiplatform viruses

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

Conditions:

$$U_1 \cap U_2 \neq 0$$

- The virus code can be the same.

# Multiplatform viruses

$$G_1 = <V_1, U_1, T_1, f_1, q_1, M_1>$$

$$G_2 = <V_2, U_2, T_2, f_2, q_2, M_2>$$

Conditions:

$$U_1 \cup U_2 \neq 0$$

- The virus code can be the same.

$$U_1 \cup U_2 = 0$$

- The virus code must be different.

# Future

# Future

- **Examine general virus detection problem in limited cases:**
  - Spreading under the model
  - Limit the time/space

# Future

- **Examine general virus detection problem in limited cases:**
  - Spreading under the model
  - Limit the time/space

- **Examine polymorphic techniques**
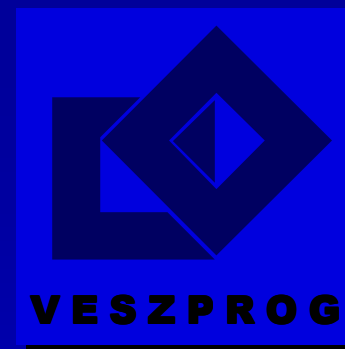  - Without coding/decoding
  - Changing instructions

# Searching technique questions

# Searching technique questions

- **For what kind of viruses can be used ?**

# Searching technique questions

- **For what kind of viruses can be used ?**

- **What is the probability of false alarms ?**

# Searching technique questions

- **For what kind of viruses can be used ?**

- **What is the probability of false alarms ?**

- **What is the expense criteria ?**

# Sequence searching algorithm

# Sequence searching algorithm

- **for non-polymorphic known viruses**

# Sequence searching algorithm

*L:* size of suspicious area
*M:* number of sequences
*N:* size of a sequence
*n:* number of values in one cell

- for non-polymorphic known viruses

- false alarms: $p \approx \dfrac{L \cdot M}{n^N}$

# Sequence searching algorithm

*L:* size of suspicious area
*M:* number of sequences
*N:* size of a sequence
*n:* number of values in one cell

- for non-polymorphic known viruses

- false alarms: $p \approx \dfrac{L \cdot M}{n^N}$

- expense criteria: P, polynomial
$$\leq L \cdot M \cdot N \text{ comparisions}$$

**VESZPROG**

# "Heuristic" algorithm

# "Heuristic" algorithm

- **for known viruses**

# "Heuristic" algorithm

- **for known viruses**

- **expense criteria:**

| Host program | Decoder  (cycle) | Body |
|---|---|---|

# "Heuristic" algorithm

- **for known viruses**

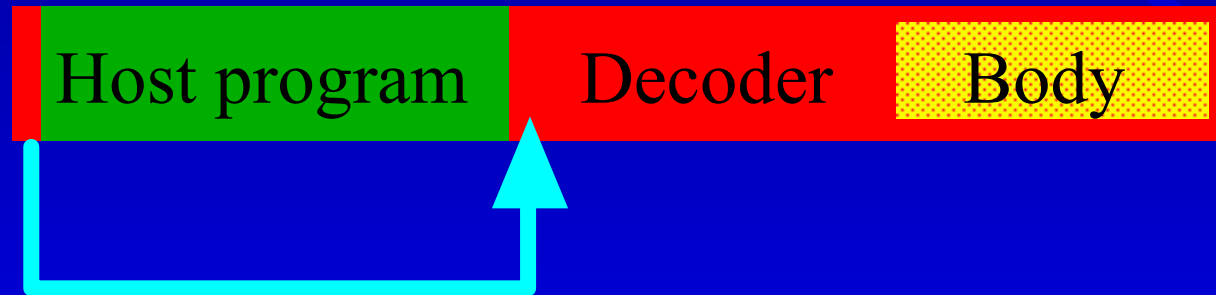- **expense criteria: NP**

$n$

| Host program | Decoder  (cycle) | Body |

Executes $2^n$ cycle !

# How can we measure the power of polymorphism ?

# How can we measure the power of polymorphism ?

Host program | Decoder | Body
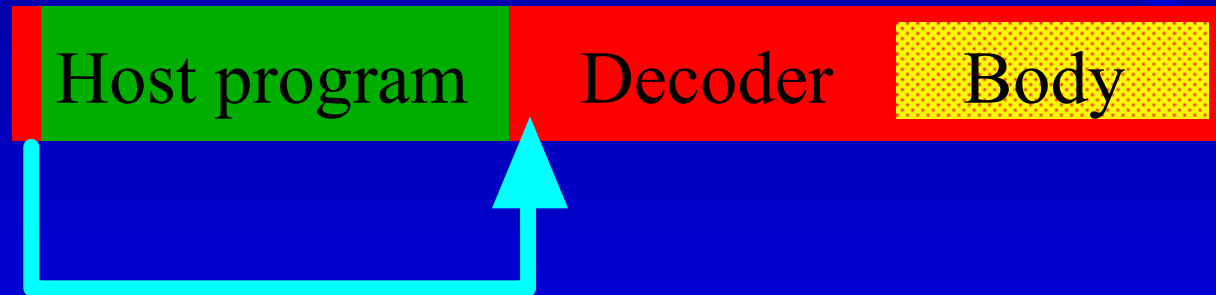
# How can we measure the power of polymorphism ?

**VESZPROG**

| Host program | Decoder | Body |
|---|---|---|

$$\alpha = \frac{\text{size of variable parts of the virus}}{\text{full size of the virus}}$$

# How can we measure the power of polymorphism ?

**VESZPROG**

| Host program | Decoder | Body |
|---|---|---|

$$\alpha = \frac{\text{size of variable parts of the virus}}{\text{full size of the virus}}$$

$$\beta = \text{number of variants of the decoders}$$

# Flowchart of a virus

# Flowchart of a virus

**VESZPROG**

search for an
uninfected program

# Flowchart of a virus

**VESZPROG**

search for an
uninfected program

append virus

# Flowchart of a virus

**VESZPROG**

search for an
uninfected program

↓

append virus

↓

choose a random
instruction in the virus

# Flowchart of a virus

**VESZPROG**

search for an
uninfected program

append virus

choose a random
instruction in the virus

swap with the next
instruction

# Flowchart of a virus

**VESZPROG**

search for an uninfected program

append virus

choose a random instruction in the virus

**repeat 100 times**

swap with the next instruction

# Flowchart of a virus

**VESZPROG**

search for an uninfected program

append virus

choose a random instruction in the **DISK**

**repeat 100 times**

swap with the next instruction

**Name:** **RIPPER**
**Aliases:** **Jack Ripper**
**Status:** **Common**
**Origin:** **Norway**
**Length:** **1024 bytes (2 sectors)**
**Infect:** **MBR, Boot sector**
**Other:** **Resident, Stealth, Disk corruption**

**VESZPROG**

# What can we do with this mathematical model ?

- **Examine the working mechanism of viruses**

- **Examine the virus detection problem**

- **Examine multiplatform viruses**

# What can we do with this mathematical model ?

**VESZPROG**

- Examine the working mechanism of viruses

- Examine the virus detection problem

- Examine multiplatform viruses

- Examine new polymorphic virus types