

# Learning Control for Batch Thermal Sterilization of Canned Foods

S. Syafie<sup>a,\*</sup> F. Tadeo<sup>a</sup> M. Villafin<sup>b</sup> A. A. Alonso<sup>b</sup>

<sup>a</sup>*Department of Systems Engineering and Automatic Control, University of Valladolid, 47011 Valladolid, Spain, {syam|fernando}@autom.uva.es*

<sup>b</sup>*Process Engineering Group IIM-CSIC, Vigo, Spain, {marcosvm|antonio}@iim.csic.es*

---

\* Corresponding author.

# Learning Control for Batch Thermal Sterilization of Canned Foods

---

## Abstract

A control technique based on Reinforcement Learning is proposed for the thermal sterilization of canned food. The proposed controller has the objective of ensuring a given degree of sterilization during Heating (by providing a minimum temperature inside the cans during a given time) and then a smooth Cooling, avoiding sudden pressure variations. For this three automatic control valves are manipulated by the controller: a valve that regulates the admission of steam during Heating, and a valve that regulate the admission of air, together with a bleeder valve, during Cooling. As dynamical models of this kind of processes are too complex and involve many uncertainties, controllers based on learning are proposed. Thus based on the control objectives and the constraints on input and output variables, the proposed controllers learn the most adequate control actions by looking up a certain matrix that contains the state-action mapping, starting from a preselected state-action space. This state-action matrix is constantly updated based on the performance obtained with the applied control actions. Experimental results at laboratory scale show the advantages of the proposed technique for this kind of processes.

### *Key words:*

Intelligent Process Control, Sterilization Process, Food Process, Batch Process, Reinforcement Learning.

---

## 1 Introduction

2 The food industries are nowadays facing critical changes in response to con-  
3 sumers, which, in addition to health and safety awareness, demand an ever  
4 larger diversity of food products with high quality standards. On the other  
5 hand, these industries are in a permanent quest for new markets and popula-  
6 tion sectors not accessible before, which immediately translates into the search  
7 for more efficient processes, in order to gain market share ([Bruin and Jongen,](#)  
8 [2003](#)).

9 This paper concentrates on the design of controllers for a specific process in  
10 the food industries, namely the so-called thermal processes for sterilization of  
11 canned foods ([Lewis, 2006](#); [Ramaswamy and Singh, 1997](#)). These processes are

12 very important for minimizing the activity of harmful microorganisms in food,  
13 thereby reducing health risks and increasing the durability of the products.  
14 For the problem at hand, the microorganism activities are reduced through  
15 thermal sterilization in pressurized retorts using steam. Unfortunately, thermal  
16 processing also produces the deterioration of the organoleptic properties of  
17 the food when conditions are not carefully controlled. For this reason, an  
18 appropriate control of the process is fundamental to guarantee the safety and  
19 quality of the products (Lewis, 2006; Ramaswamy and Singh, 1997).

20 Thus, the central objective of controllers for the sterilization process is the in-  
21 activation of microorganisms present in the foodstuff, while preserving as much  
22 as possible product quality, avoiding very quick variations in temperature and  
23 pressure and minimizing the operation time. For this, the sterilization process  
24 can be divided in three stages that use different control strategies: Venting,  
25 Heating and Cooling. Venting is normally carried out manually, so the stages  
26 of the process relevant from the point of view of controller design are Heat-  
27 ing (where the main objective is to ensure a given degree of sterilization by  
28 ensuring a given temperature during a certain time by manipulating the en-  
29 trance of steam in the retort), and Cooling (where the temperature is carefully  
30 decreased by replacing the steam with air).

31 The kinetics of thermal destruction of microorganisms or degradation of nu-  
32 trients are usually assumed to follow pseudo-first-order kinetics (e.g. the TDT  
33 model) with an exponential-type temperature dependence (Balsa-Canto et al.,  
34 2002a,b). Such kinetics constitutes the basis to quantify the degree of steril-  
35 ization, usually given in terms of *lethality* (in units of time), that defines the  
36 amount of time required to produce a certain decimal reduction. For details,  
37 the reader is referred to Ramaswamy and Singh (1997). Unfortunately, due  
38 to the complexity of the process, the variability of the products to be ster-  
39 ilized and the reduced number of sensors it is not feasible to derive models  
40 adequate for model-based controller design. To deal with this issue, this pa-  
41 per concentrates on the application of a control technique based on learning.  
42 More precisely, a Model-Free Learning Controller (MFLC) will be developed for  
43 this thermal sterilization processes. This MFLC is based on *Reinforcement*  
44 *Learning*, so it is an agent-based technique based on re-framing the problem  
45 of achieving process control objectives by learning through interaction with  
46 the process (see Figure 1), taking always into account the inherent constraints  
47 in input and output signals. The (*agent*) interacts with the rest of the process  
48 (also called *environment* in learning approaches): the agent selecting actions  
49 and the environment responding to those actions and presenting new situations  
50 to the agent. The environment also provides rewards, that are numerical values  
51 that the agent tries to maximize, as they give a measurement of performance  
52 (Sutton and Barto, 1998). More specifically, the agent and the environment  
53 interact at each of a sequence of discrete time step. At each time step, the  
54 agent receives some representation of the environment's state, and on that

55 basis selects an action. The agent receives a numerical reward, and moves to  
56 a new state (Sutton and Barto, 1998). Thus, the reward function depends on  
57 the recent state, action and successor state: with time, the agent gathers more  
58 information and provides optimal actions for every visiting state.

59 Although Reinforcement Learning ideas seem promising, they were not de-  
60 veloped for process control problems (Sutton and Barto, 1998; Bertsekas and  
61 Tsitsiklis, 1996), so in this paper the Model-Free Learning Control (MFLC)  
62 technique (Syafie et al., 2007a; Syafie et al., 2007b) is used to control the  
63 sterilization process. This MFLC is gives a feasible implementation of Rein-  
64 forcement Learning for process control problems, by providing a precise but  
65 simple definition of symbolic states and actions, based on control objectives  
66 and the constraints on input and output variables. This methodology is com-  
67 plementary to other intelligent control approaches (such as Fuzzy Logic or  
68 Neural Networks), in the sense that initial values for the parameters of the  
69 MFLC algorithm can be derived from previous controllers. Starting from these  
70 initial parameters, using learning MFLC provides a simple methodology to im-  
71 prove the controller by interaction with the plant.

72 The rest of this article is structured as follows: First the background and scope  
73 are stated in Section 2. A short presentation of the thermal sterilization pro-  
74 cess is given in Section 3. The proposed technique to control the sterilization  
75 process by using Model-Free Learning Control (MFLC) is given in Section  
76 4. The MFLC application for controlling a sterilization process at laboratory  
77 scale is discussed in Section 5. Finally, some conclusions are given in Section  
78 6.

## 79 2 Background and scope

80 In industrial sterilization processes for canned food the most common con-  
81 trollers are still PID. For example in Mulvaney et al. (1990), a Proportional  
82 Integral (PI) controller was developed for this process. A study using a com-  
83 bination of the linearizing-transformation of differential geometry and the  
84 quality-control of Q-PID/Q-PI was presented by Alonso et al. (1993), whereas  
85 a PID-type controller with parameters selected using Internal Model Control  
86 (IMC) was reported by Alonso et al. (1997, 1998). It was found that PID  
87 controllers work well during Heating as long as the plant is operated in small  
88 neighborhoods of the constant-heating temperature around the tuning region;  
89 unfortunately, frequently the controllers have to be retuned to operate in other  
90 conditions (for example, when the type and amount of cans change) , which  
91 is cumbersome.

92 Advanced control strategies have also been proposed for this process, such as

93 the online correction of the lethality value reported by [Teixeira and Tucker](#)  
94 [\(1997\)](#). In [Kuma et al. \(2001\)](#), an algorithm based on three control modes was  
95 presented, but no specific proposal was given on how to regulate the steam,  
96 water, drain, air and bleeder valves. An optimal control problem with state  
97 and control constraints governed by a nonlinear heat equation was proposed  
98 by [Kleis and Sachs \(1999\)](#). The discretized optimal control was expressed as  
99 a large-scale continuous optimization, which can be solved using sequential  
100 quadratic programming. However, the proposed algorithm was mathemati-  
101 cally complicated. A closed-loop optimal receding horizon controller (RHC)  
102 incorporating model uncertainty was designed and studied by [Chalabi et al.](#)  
103 [\(1999\)](#), where a non-gradient method was used to solve the corresponding non-  
104 linear optimization problem. Unfortunately, this kind of controllers requires  
105 that all the states of the system to be measurable, which is impractical. Since  
106 all these advanced controllers are difficult to design and need a precise mathe-  
107 matical model of the process, the most frequent control technique in industry  
108 is still, therefore, a manual supervision of PID controllers.

109 To deal with problems of batch to batch variations and the complexity of  
110 the models for control, techniques based on learning would be adequate as  
111 they adapt to the specific situation at hand through the result of previous  
112 experiences. Techniques based on Reinforcement Learning have been selected,  
113 as they provide a rigorous methodology for learning without detailed mathe-  
114 matical models of the controlled plant, using a simple algorithm suitable for  
115 real-time implementation ([Sutton and Barto, 1998](#)).

116 In particular the MFLC approach, previously proposed by some of the authors  
117 ([Syafie et al., 2007a](#); [Syafie et al., 2007b](#)), will be used to control the thermal  
118 processing, as it corresponds to a feasible implementation of Reinforcement  
119 Learning algorithms ([Sutton and Barto, 1998](#)) for Process Control. This tech-  
120 nique is used because it is simple and does not need a precise *a priori* model of  
121 the process, but incorporates basic knowledge of the process behavior (infor-  
122 mation from output range, control limitations, loop interactions, etc). Thus,  
123 in MFLC controllers the control objective is expressed as the optimization of  
124 a desired performance index by learning to apply appropriate control actions  
125 through interaction with the plant. In particular, the MFLC approach pro-  
126 posed here is based on  $Q$ -learning ([Sutton and Barto, 1998](#); [Bertsekas and](#)  
127 [Tsitsiklis, 1996](#)). However, the idea can be easily augmented to improve learn-  
128 ing speed by applying other methodologies in literature, such as lazy learning  
129 ([Atkenson et al., 1997a,b](#)), near optimal closed-loop control ([Ernst, 2003](#)) and  
130  $q$ -iteration with CMACS ([Timmer and Riedmiller, 2007](#)).

131 We must point out that, although for simplicity, and in order to represent  
132 industrial practice, the problem at hand is represented as a sequence of two  
133 dynamical systems (during Heating a single-input single-output system, and  
134 during Cooling a two-input single-output system), if needed the proposed ap-

135 proach can be extended to more complex multiple-input multiple-output sys-  
136 tems using the ideas of Riedmiller (1997).

### 137 3 Batch Thermal Sterilization Process

138 The thermal sterilization processes for prepackaged food can be carried out in  
139 continuous or batch units. This article concentrates on learning to control the  
140 thermal sterilization process in batch units, as it is the most frequent approach  
141 in the industry, and the one that can make better use of a learning approach.

#### 142 3.1 Process Description

143 The sterilization process is assumed to be carried out in batch steam retorts as  
144 depicted in Figure 2. A typical operation cycle involves several stages, which  
145 in this paper are assumed to be the following:

- 146 • *Venting*: In this initial stage, steam is introduced in the retort to eliminate  
147 the air, so heat transmission is more efficient during Heating. At this stage,  
148 bleeder and drain valves are open. When the pressure in the retort,  $P_r$ ,  
149 matches that corresponding to saturated steam,  $P_s$ , at that temperature,  
150 there is only steam in the retort, so Heating can start.
- 151 • *Heating*: The objective of this central stage is that the temperature inside  
152 the retort is at the level required, for enough time to reach the desired  
153 microbiological lethality. At time  $t$  the lethality  $F(t)$  is defined as follows:

$$154 \quad F(t) = \int_0^t 10^{\frac{T(t)-T_{ref}}{z_{ref}}} dt \quad (1)$$

155 where  $z_{ref}$  and  $T_{ref}$  are parameters that depend on the container and the  
156 product, which are obtained experimentally, and  $T(t)$  is the temperature at  
157 the critical point (the point inside the product with lowest temperature),  
158 (see Ramaswamy and Singh (1997); Alonso et al. (1997)). This lethality  
159 is affected by small variations in the temperature, so automatic control is  
160 required during this cycle.

- 161 • *Cooling*: Once the Heating period concludes, the product is cooled with  
162 water down to room temperature. At the same time, air is injected into  
163 the retort to avoid sudden pressure drops that could result in the bursting  
164 of the product containers. Pressure control during this stage is especially  
165 important for glass containers or conduction heated-type products where the  
166 existence of sharp temperature gradients between the inside and the outside  
167 of the product induces high differential pressure (Alonso et al., 1997, 1998).

## 168 4 MFLC Technique

169 The Model-Free Learning Control technique (MFLC) that is proposed here for  
170 batch sterilization processes is a control technique, based on Reinforcement  
171 Learning (Sutton and Barto, 1998; Bertsekas and Tsitsiklis, 1996), which gives  
172 a feasible implementation of automatic learning in process control problems,  
173 by providing a precise definition of symbolic states and actions, based on  
174 control objectives and the constraints on input and output variables. It has  
175 been presented in detail by the some of the authors in Syafie et al. (2007a);  
176 Syafie et al. (2007b), so only the main ideas are given here.

### 177 4.1 MFLC Architecture

178 The MFLC architecture is represented in Figure 3: as with most Reinforcement  
179 Learning algorithms, it is based on describing the system in terms of symbolic  
180 states, so the controller learns how good the application of a given action in  
181 a given state is, by applying the action to the system and then checking the  
182 quality of the response. The evaluation of the effect of each action is done  
183 by estimating the expected return mathematically, storing the values of this  
184 return (which measure the quality of the response) in the so-called  $Q$ -matrix  
185 (discussed in section 4.2).

186 The MFLC is based on a precise selection of states, actions and control signals  
187 (discussed in sections 4.3 and 4.4), with the objective of representing typical  
188 problems in process control and being easily understood by the final user.  
189 The operation of the algorithm, represented in Figure 3 is based on, first, the  
190 selection of the agent of one action from those available in the current state,  
191 using the "Policy". Then, the action is converted to a control signal in the  
192 "Calculation U" block. Then, based on the measured output, the "Situation"  
193 block estimates the next state and the corresponding reward. From this re-  
194 ward, the so-called  $Q$ -value is updated in the "Critic" block, which reflects  
195 the adequacy of the action applied. As time goes by, actions are selected by  
196 the agent, and learning is carried out by checking the quality of the response:  
197 Actions that drive the system into the goal state are considered to be good,  
198 so its  $Q$ -value is increased. On the other hand, actions that do not drive the  
199 system into the goal state are punished.

### 200 4.2 $Q$ -matrix

201 Mathematically, the objective in MFLC is to maximize the expected return  
202 (Sutton and Barto, 1998) taking into account the control and state constraints.

203 A central part of the learning algorithm is the estimation of this expected  
 204 return. For this, the state-action value function,  $Q(s, a)$ , is used, as it contains  
 205 the expected return, when starting from the state  $s$ , the agent applies the  
 206 action  $a$ , and thereafter follows the policy  $\pi$ :

$$207 \quad Q^\pi(s, a) = E_\pi\{R_t | s_t = s, a_t = a\} = E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right\}. \quad (2)$$

208 This function is stored in a matrix  $Q(s_t, a_t)$ , the  $Q$ -matrix. At each sampling  
 209 time, these  $Q$ -values are calculated by taking into account the current and  
 210 future benefits: when action  $a_t$  has been selected and applied to the plant, the  
 211 system moves to a new state,  $s_{t+1}$ , and receives a reinforcement signal,  $r_{t+1}$   
 212 (which evaluates the quality of the response), so the  $Q$ -matrix is updated as  
 213 follows:

$$214 \quad Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_{b \in A_{s_{t+1}}} Q(s_{t+1}, b)] \quad (3)$$

215 where:

- 216 - The learning rate,  $\alpha \in (0, 1]$ , is a tuning parameter that can be used to  
 217 optimize the speed of learning (a large learning rate makes learning faster,  
 218 but might induce oscillations). It is required for computation of expectation  
 219 in the form of an iterative averaging.
- 220 - The discount factor,  $\gamma \in (0, 1]$ , is used as a factor to weight the effect more  
 221 heavily in the near future: If  $\gamma$  is small, the agent learns to behave only for  
 222 short-term reward; the closer  $\gamma$  is to 1 the greater the weight assigned to  
 223 long-term reinforcements.
- 224 -  $A_{s_{t+1}}$  is the finite set of possible actions in the new state.

### 225 4.3 State Representation

226 A central issue in all Reinforcement Learning algorithms is the definition of the  
 227 states, which are symbolic and represent the "distance" to the goal. In MFLLC,  
 228 the states are defined based on the control objective and the constraints on the  
 229 control signal and the states, as follows: the control objective is considered to  
 230 be to maintain the desired output inside the band  $r - d$  and  $r + d$ , as shown in  
 231 Figure 4. The width of this band is defined based on the tolerance of the system  
 232 (which depends on measurement noise, disturbances and the specifications).  
 233 This band is defined as the *goal band*, and corresponds to the *goal state*, where  
 234 the agent should drive the system and ensures that it remains there (it is  
 235 now assumed, without loss of generality, that is exactly in the middle of the  
 236 working range). To describe the rest of the symbolic states, it is considered



237 that the agent has  $h$  states from the goal state to the maximum positive  
 238 or minimum negative error of the system,  $f$  (Selecting  $h$  is a trade-off: this  
 239 number must be large enough to describe all the different responses of the  
 240 process, but small enough to reduce computational time and the size of the  
 241  $Q$ -matrix). The "span" of each state can be calculated as follows:

$$242 \quad c = \frac{f - d}{h}. \quad (4)$$

243 Thus, the positive bound parameter can be presented as:

$$244 \quad \omega_i = d + (i - 1)c, i \in [1, \dots, h] \quad (5)$$

245 (For negative errors, the bound parameter is trivial by changing signs). Thus,  
 246 the vector of symbolic states can be presented as follows:

$$247 \quad g_j = \begin{cases} e - \omega_j & \text{if } e \leq \omega_j; \\ \omega_j - e & \text{else,} \end{cases} \quad j \in [1, \dots, 2h + 1] \quad (6)$$

248 where  $e$  is the tracking error. The symbolic current state,  $s_t$ , is just:

$$249 \quad s_t = \arg \max_j (g_j). \quad (7)$$

#### 250 4.4 Action Representation

251 In the single-input single-output version of MFLC, the control signal  $u_t \in \mathbb{R}$   
 252 is calculated by varying the previous control signal in a magnitude calculated  
 253 from the difference of the numerical values of the selected optimal action,  $a_t \in$   
 254  $\mathbb{N}$ , with respect to the *wait action*,  $a_w$  (action corresponding to maintaining  
 255 the previous control signal). That is:

$$256 \quad u_t = u_{t-1} + k(a_w - a_t), \quad (8)$$

257 where  $k$  is the tuning parameter. This gives a PI-like structure, which simplifies  
 258 initialization and tuning for the end user. At each state there is only a finite  
 259 set of possible actions (see Figure 5). These actions are selected based on the  
 260 systems description: in particular, from the limitations on the minimum and  
 261 maximum variations of the control signal, as follows: Let the control variations  
 262 be bounded as follows:

$$263 \quad \underline{\Delta u} \leq |\Delta u| \leq \overline{\Delta u}, \quad (9)$$

264 where  $\underline{\Delta u}$  and  $\overline{\Delta u}$  are known bounds. The number of total actions needed to  
 265 satisfy the constraints can be calculated as follows:

$$266 \quad N_a = 2h \left( \text{round} \left( \frac{\overline{\Delta u} - \underline{\Delta u}}{kh} \right) \right) + 1, \quad (10)$$

267 where the round-up function is used. From (8), (9) and (10), the value corre-  
 268 sponding to the wait action  $a_w$ , can be calculated as follows:

$$269 \quad a_w = \frac{N_a + 1}{2}. \quad (11)$$

270 If there is no overlapping, the number of actions in each state can be calculated  
 271 being  $n_a = \frac{N_a - 1}{2h}$ . However, to increase the number of available actions and  
 272 represent nonlinear action-to-space relations (important in process control),  
 273 a degree of overlapping must be included (see Figure 5). Of course, at each  
 274 state, not all the actions are available: Each state has a subset of actions. For  
 275 example, during Heating, if the tracking error for temperature is very small,  
 276 the only actions available are those that increase to correct the temperature.  
 277 Thus, the number of actions in each state is  $n_a^\beta = n_a(1 + \beta)$ , where  $\beta$  is a  
 278 parameter that gives the degree of overlapping with neighboring states (always  
 279 selected such that  $n_a^\beta$  is integer). Then, the available actions for every state go  
 280 from  $a_p^j$  to  $a_b^j$  (except in the goal state, where there is only the wait action).  
 281 The idea is presented in Figure 5 and developed in [Syafie et al. \(2008\)](#). Those  
 282 available actions can be calculated as

$$283 \quad \begin{aligned} a_p^j &= a_p^{j-1} + (j - 1)v, \\ a_b^j &= a_p^j + n_a^\beta - 1, \end{aligned} \quad (12)$$

284 where  $v = \beta \frac{n_a^\beta}{h}$  and  $a_p^{j-1}$  is the first action in the state  $j$  calculated as

$$285 \quad a_p^{j-1} = \begin{cases} 1, & \text{if } j = 1 \\ 2a_w - a_b^{j-2}, & \text{if } j = h + 2 \end{cases}. \quad (13)$$

286 The strategy for selecting one action from those available ones is through  
 287 *exploration* and *exploitation* policies. The agent explores those available actions  
 288 to know the optimal value function by executing trial actions, following the  
 289  $\varepsilon$ -greedy policy ([Sutton and Barto, 1998](#)). This means that the action which  
 290 has the maximum  $Q$ -value will be selected with  $1 - \varepsilon$  probability and the rest  
 291 will explore trial actions selected from those available in the state.

## 292 5 Thermal Control of Prepackaged Food

293 This section explains the application of MFLC ideas for batch thermal ster-  
294 ilization. The first part of this section discusses the control strategy, followed  
295 by a discussion on the selection of the parameters of the controllers for the  
296 Heating and Cooling stages of these sterilization processes.

297 As discussed in Section 3, there are three crucial steps in controlling the ster-  
298 ilization process: Venting, Heating and Cooling.

299 The proposed control strategy for these cycles is shown in Figure 6. As the  
300 venting stage can be controlled using a simple technique (keeping bleeder and  
301 drain valves fully open until the pressure inside the retort  $P_r$  reaches the  
302 steam pressure  $P_s$ ), the control application therefore concentrates on Heating  
303 and Cooling. The use of MFLC for Heating and Cooling is now presented.

### 304 5.1 Heating Control Strategy

305 During Heating, the control objective is to maintain the temperature inside  
306 the goal band by manipulating the steam valve. To evacuate the condensed  
307 water from the retort, the drain valve is open. Also, the bleeder valve is slightly  
308 open.

309 Mathematically, during Heating, the objective is to maintain the retort tem-  
310 perature within a tolerance of  $\pm 2.0^\circ\text{C}$  with respect to the provided reference.  
311 Thus, the goal band is  $r - 2.0$  to  $r + 2.0$ . The output range is considered  
312 to be  $\pm 4.0^\circ\text{C}$  with respect to the reference. Thus, from these numbers and  
313 following the ideas presented in Section 4, there are 21 symbolic states, where  
314 state #11 corresponds to the goal state. The actions are then defined based  
315 on the possible control variations: the signal must vary within the following  
316 bounds:

$$317 \quad 0.0001 \leq |\Delta u| \leq 0.008. \quad (14)$$

318 Thus, the  $Q$  matrix size is  $1601 \times 21$ , where the wait action is action #801 (this  
319 matrix will be denoted  $Q_H$ ). The tuning parameter is selected to be  $k = 10^{-5}$ ,  
320 based on the control constraints. To include some nonlinearity, a small overlap  
321 is considered, with the number of actions in every symbolic state to be 158.  
322 Therefore, in state #1 the actions are #1,  $\dots$ , #158, in state #2 the actions  
323 are #71,  $\dots$ , #228, and so on, following (12). The controller parameters are  
324 summarized in Table 1.

325 The objective of the control task is to maintain the process in the goal state,  
 326 or return it to the goal state if there has been any disturbance or change of  
 327 reference. To achieve this, maximum reward is introduced for actions causing  
 328 the process error to be smaller than the previous one. Actions that move the  
 329 system away from the goal band are punished. Therefore, the reward is given  
 330 as:

$$331 \quad R_t = \begin{cases} 1.0 & \text{if } |e_t| \leq |e_{t-1}|, \\ -1.0 & \text{otherwise.} \end{cases} \quad (15)$$

332 Of course, more complex reward functions could be selected, but this particular  
 333 reward function has been selected following the ideas in [Smart \(2002\)](#), which  
 334 recommends not indicating a detailed path for the agent to achieve the goal,  
 335 but only the goal, as the path assumed to be the most adequate might not  
 336 really be the best (learning takes care of finding the most adequate approach).  
 337 Thus, this gives an approach parallel to the Mayer-type objective functions in  
 338 Optimal Control ([Stryk and Bulirsch \(1992\)](#)), with the trajectory constrained  
 339 by the limited number of actions available in each state.

340 Heating finishes when the desired lethality time  $t_l$  is reached (where  $t_l$  is  
 341 evaluated from (1)). That is, denoting by  $t_v$  the starting time of the Heating,  
 342 the agent switches from Heating control to Cooling control when  $t \geq t_v + t_l$ .

## 343 5.2 Cooling Control Strategy

344 The state-action space has been discussed in detail for the Heating stage in  
 345 Section 5.1. In the Cooling stage, the objective of the controller design is  
 346 to avoid sudden pressure drops by regulating air and bleeder valves. The air  
 347 valve is used to increase or maintain pressure, while the bleeder valve is used to  
 348 reduce the pressure inside the retort. Avoiding sudden pressure drops is aimed  
 349 at avoiding food container bursts. On the other hand, the food containers are  
 350 cooled down to room temperature. This is achieved by flowing water into  
 351 the retort. In this stage, the water stream is set with a fixed stream. When  
 352 the retort temperature is reached, the water flow is cut off. To avoid large  
 353 disturbances at the beginning of the Cooling stage, the steam present in the  
 354 retort is gradually eliminated. However, the drain valve is kept open.

355 To select the structure of the  $Q$ -matrix for this stage, denoted now  $Q_C$ , a  
 356 similar strategy as in Section 5.1 is used. Since there are two control signals  
 357 (the Air and Bleeder valves), this  $Q_C$ -matrix is designed with three dimensions  
 358 (one state for each combination of two actions): The matrix represents the  
 359 space of error in the pressure to the air-valve-action and the bleeder-valve-

360 action.

361 The control parameters for the Cooling state are shown in Table 2. Even  
362 though the same controller gain,  $k$ , is used in the design of the air and bleeder  
363 action spaces, the gain, can, however be tuned separately in implementation.

## 364 6 Results and Discussion

365 This section discusses the application of the proposed MFLC controller for  
366 controlling thermal canned food sterilization in a laboratory plant, placed at  
367 the Maritime Research Center, Vigo, Spain. The agent-based MFLC is initial-  
368 ized by training using a virtual plant (simulation). Then, online application is  
369 implemented at the laboratory-scale autoclave.

### 370 6.1 Plant Description

371 A schematic of the batch retort unit used for testing the algorithms developed  
372 in this paper is presented in Figure 2. The vessel, built in steel, has an approx-  
373 imate weight of 150 kg, and dimensions of approximately 1m of length and 60  
374 cm of diameter. To record the evolution of the relevant variables during pro-  
375 cessing, three PT100, eight thermocouples and a pressure sensor are located  
376 inside the vessel. A computer system is used to gather and analyse real time  
377 data. Process Control is carried out using Labview, with an external module  
378 WebDAQ that connects the PT100 and pressure sensors to the controller by  
379 means of an Ethernet port, and an ADAM that connect the thermocouples. A  
380 NiDAQ card is used to actuate the valves, that are Siemens PV90 (DN15)-flat  
381 seat, with nominal linear characteristics.

### 382 6.2 Initial Training of the Agent

383 The detailed model of the thermal canned-food process using a retort proposed  
384 in [Alonso et al. \(1997\)](#) was used to train the  $Q_H$  and  $Q_C$  matrices. The model,  
385 based on nonlinear dynamic equations, was numerically written and solved  
386 in Ecosimpro<sup>®</sup> simulation language ([Ecosimpro, 1999](#)), with training done  
387 for various learning stages. The main reasons for using a virtual plant for  
388 initial training are the reduction of costs and the prevention of damage to  
389 the products during learning for extreme situations. If a simulation were not  
390 available, the  $Q_H$  and  $Q_C$  matrices can be initialized adapting values from  
391 similar processes or using values from previous controllers.

392 The temperature and pressure responses of the first training stage using the  
393  $Q_H$ -matrix are shown in Figures 7 and 8. During Heating, the control objective  
394 is to maintain a given pre-selected time-temperature profile so as to ensure the  
395 appropriate lethality by manipulating the steam valve. Note that the pressure  
396 does not need to be controlled during this stage, since the steam is saturated  
397 and no air is present in the retort after venting.

398 After the lethality time  $t_l$  is satisfied, the system enters the Cooling stage.  
399 In this stage, the temperature is not controlled. In other words, there is no  
400 valve regulation rule for controlling the temperature. So that the canned food  
401 reaches a cool temperature (approximately ambient temperature), water is  
402 passed into the retort at a fixed rate. The water valve is then gradually opened  
403 up to 30%. The valve opening in this position is to avoid flooding inside the  
404 retort and to provide enough water for cooling. In this Cooling stage, the  
405 objective of the controller is switched to control the pressure (see Figure 7b).  
406 To avoid sudden pressure drops, the air valve is initially fully open. At the  
407 same time, the bleeder valve is totally closed, to avoid losing air inside the  
408 retort. Both air and bleeder valves are regulated according to the pressure  
409 measured inside the retort. The last pressure reading of the Heating stage is  
410 used as an initial pressure set point. From this initial reference, the pressure  
411 reference is gradually reduced by 500 Pa if the system is inside the goal state  
412 and/or above  $10^5$  Pa. This value can be changed according to the resistance  
413 of the container material. After some training stages, the  $Q_C$ -matrix is used  
414 in the online implementation.

415 The agent is also trained for some environment changes, such as changes in  
416 the temperature of reference (Figure 9). The learning control is able to track  
417 the set point changes and correct the error. Finally, the responses are inside  
418 the desired region.

### 419 6.3 Application on the laboratory process

420 The online implementation of MFLC for controlling temperature and pressure  
421 of the canned food process is discussed in this section. As mentioned above,  
422 the feedback signals are the average temperature in the basket and the average  
423 pressure. Temperature responses during the Heating stage are shown in Figure  
424 10a, and the pressures inside the retort are plotted in Figure 10b. The control  
425 signal is depicted in Figure 11: only the steam valve position is plotted, as the  
426 other valves remain constant. It can be seen that the steam valve works within  
427 the range from 0 to 20% opening. Therefore, the control signal is bounded  
428 within the desired range. In this application, the steam flow is equipped with  
429 a relief valve to reduce the pressure. Hence, the maximum pressure of the  
430 steam entering the retort is always about 2 atm.

431 In summary, adequate temperature control for the Heating process was ob-  
432 tained in the laboratory plant. From the laboratory application, the proposed  
433 learning control is able to track the temperature and keep it inside the desired  
434 bound (Figure 10 a) during the Heating stage. Also, the controller is able  
435 to regulate the system for setpoint changes, while the temperature remains  
436 within the desired bounds. The controller output for the setpoint regulation is  
437 presented in Figure 11. The controller manipulates the steam valve smoothly,  
438 with a control signal suitable for the regulation of the motorized valves.

439 After a relatively short time (approximately 7 minutes for settling time), the  
440 controller can bring the system to be and remain inside the desired bound,  
441 with only a small overshoot. The performances of the proposed controller are  
442 summarized in Table 3.

## 443 **7 Conclusions**

444 A procedure for automatic control of the sterilization process in canned food  
445 industry has been presented, based on the use of controllers based on learning.  
446 More precisely, a controller is proposed to manipulate the steam valve during  
447 Heating, using the Model-Free Learning Control (MFLC) strategy, followed by  
448 another MFLC controller to regulate the air and drain valves during Cooling.  
449 The results of the application of the methodology in a plant at laboratory  
450 scale show that the proposed controllers make it possible to maintain the  
451 temperature and pressure of the sterilization process within specifications,  
452 allowing the safe consumption of the food.

## 453 **Acknowledgement**

454 This work has been funded by project DPI2007-66718-C04-02. The authors  
455 are thankful to Prof. Ernesto Martinez, Dr. Carlos Vilas and Dr. Miriam  
456 Rodriguez for many helpful discussions.

## 457 **References**

458 A. A. Alonso, R. I. Perez-Martin, N. V. Shukla, and P. B. Deshpande, On-  
459 line quality control of non-linear batch systems: application to the thermal  
460 processing of canned foods, *Journal of Food Engineering*, Vol. 19, pp. 275-  
461 289, 1993.

- 462 A. A. Alonso, J. R. Banga and R. P. Martin, A Complete Dynamic Model for  
463 the Thermal Processing of Bioproducts in Batch Units and its Application  
464 to Controller Design, *Chemical Engineering Science*, vol. 52, no. 8, 1997,  
465 pp. 1307-1322.
- 466 A. A. Alonso, J. R. Banga and R. P. Martin, Modeling and Adaptive Control  
467 for Batch Sterilization, *Computers and Chemical Engineering*, vol. 22, no.  
468 3, 1998, pp. 445-458.
- 469 C. G. Atkenson, A. W. Moore and S. Schaal, Locally Weighted Learning,  
470 *Artificial Intelligence Review*, Vol. 11, pp. 11-73, 1997a
- 471 C. G. Atkenson, A. W. Moore and S. Schaal, Locally Weighted Learning for  
472 Control, *Artificial Intelligence Review*, Vol. 11, pp. 75-113, 1997b
- 473 E. Balsa-Canto, A. A. Alonso, J. R. Banga, A novel, efficient and reliable  
474 method for thermal process design and optimization. Part I: theory, *Journal*  
475 *of Food Engineering*, Vol. 52, pp. 227 -234, 2002a.
- 476 E. Balsa-Canto, A. A. Alonso, J. R. Banga, A novel, efficient and reliable  
477 method for thermal process design and optimization. Part II: Application,  
478 *Journal of Food Engineering*, Vol. 52, pp. 235 -247, 2002b.
- 479 D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming*, Athena  
480 Scientific, Belmont, Massachusetts, 1996.
- 481 S. Bruin, Th. R. G. Jongen, Food Process Engineering: the last 25 years and  
482 challenges ahead, *Comprehensive Reviews in Food Science and Food Safety*,  
483 Vol. 2, 2003, pp. 42-81.
- 484 Z. S. Chalabi, L. G. van Willigenburg and G. van Straten, Robust Optimal  
485 Receding Horizon Control of the Thermal Sterilization of Canned Foods,  
486 *Journal of Food Engineering*, vol. 40, 1999, pp. 207-218.
- 487 EA Int (1999), *EcosimPro User Manual*, Available from: [www.ecosimpro.com](http://www.ecosimpro.com)  
488 (accessed 17 June 2010).
- 489 Damien Ernst, *Near optimal closed-loop control. Application to electric power*  
490 *systems*, PhD thesis at University of Liège, Belgium, 2003.
- 491 D. Holdsworth and R. Simpson, *Thermal Processing of Packaged Foods*,  
492 Springer Verlag, 2007.
- 493 D. Kleis and E. W. Sachs, Optimal Control of the Sterilization of Prepackaged  
494 Food, *SIAM Journal on Optimization*, vol. 10, no. 4, 1999, pp. 1180 - 1195.
- 495 M. A. Kumar, M. N. Ramesh and S. Nagaraja Rao, Retrofitting of a Vertical  
496 Retort for On-line Control of the Sterilization Process, *Journal of Food*  
497 *Engineering*, vol. 47, 2001, pp. 89 - 96.
- 498 M. J. Lewis, Thermal Processing, in J.G. Brennan, *Food Processing Handbook*,  
499 Wiley, 2006, pp. 33-70.
- 500 S. J. Mulvaney, S. S. H. Rizvi, and C. R. Johnson, Dynamic modelling and  
501 computer control of a retort for thermal processing, *Journal of Food Engi-*  
502 *neering*, Vol. 11, pp. 273 - 289, 1990
- 503 H. S. Ramaswamy and R. P. Singh, Sterilization Process Engineering, in K. J.  
504 Valentas, E. Rotstein, R. P. Singh, *Handbook of Food Engineering Practice*,  
505 CRC Press, New York, 1997.
- 506 Martin Riedmiller, Learning Control for Continuous MIMO Dynamical Sys-



- 507       tems Using Neuro Dynamic Programming, in proceeding of *Third European*  
508       *Workshop on Reinforcement Learning*, Rennes, France, October 13-14, 1997.
- 509 William Donald Smart, *Making Reinforcement Learning Work on Real Robots*,  
510       PhD thesis at Brown University, 2002.
- 511 O. von Stryk and R. Bulirsch, Direct and Indirect Methods for Trajectory  
512       Optimization, *Annals of Operations Research*, vol. 37, 1992, pp. 357-373.
- 513 R. S. Sutton and A. G. Barto, *Reinforcement Learning: an Introduction*, The  
514       MIT Press, Cambridge, MA, 1998.
- 515 S. Syafie, F. Tadeo and E. Martinez, Learning to Control pH Processes at Mul-  
516       tiple Time Scales: Performance Assessment in a Laboratory Plant, *Chemical*  
517       *Product and Process Modeling*, vol. 2, no. 1, 2007a, article no 7.
- 518 S. Syafie, F. Tadeo and E. Martinez, Model-Free Learning Control of Neutral-  
519       ization Process Using Reinforcement Learning, *Engineering Applications of*  
520       *Artificial Intelligence*, Vol. 20, pp. 767 – 782, 2007b.
- 521 S. Syafie, Garcia M., Vilas C., Alonso A., Martinez E., Tadeo F., Intelligent  
522       Control Based on Reinforcement Learning of Batch Thermal Steriliza-  
523       tion of Canned Foods, in *Proceedings of the 17th IFAC World Congress*, Seoul,  
524       South Korea, July 2008.
- 525 A. A. Teixeira and G. S. Tucker, On-line Retorts Control in Thermal Steril-  
526       ization of Canned Foods, *Food Control*, vol. 8, no. 1, 1997, pp. 13 - 20.
- 527 S. Timmer and M. Riedmiller, Fitted Q Iteration with CMACs, in *Proceedings*  
528       *of the International Symposium on Approximate Dynamic Programming and*  
529       *Reinforcement Learning (ADPRL)*, Honolulu, USA, April 2007.

Fig. 1. Agent-environment interaction

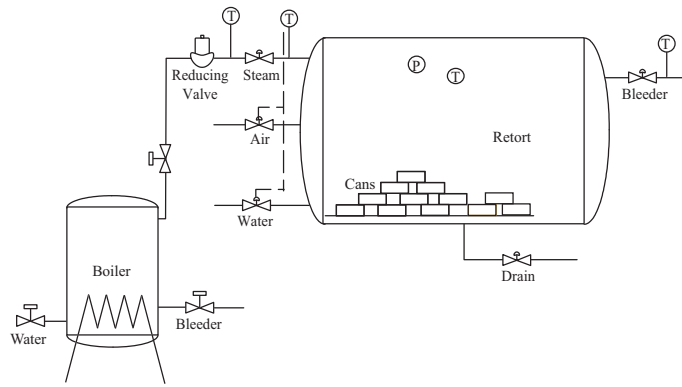


Fig. 2. Schematic of batch sterilization for controller design

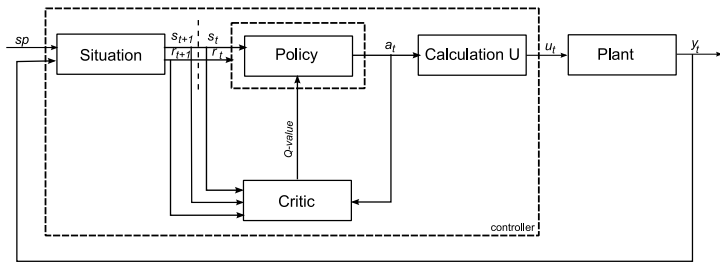


Fig. 3. MFLC architecture

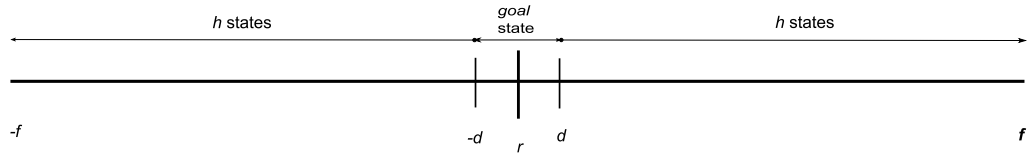


Fig. 4. Definition of the symbolic states in MFLC

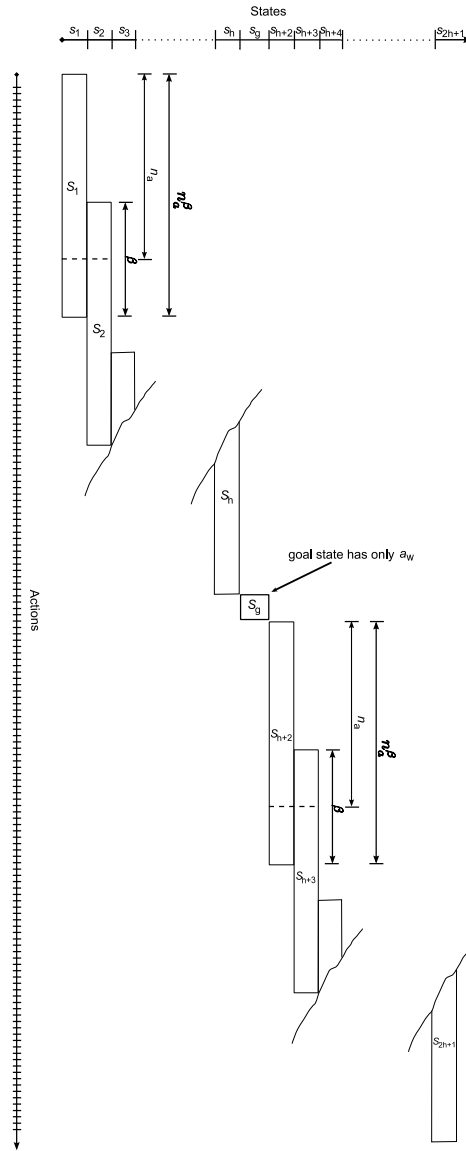


Fig. 5. State-Action space of  $Q$ -matrix

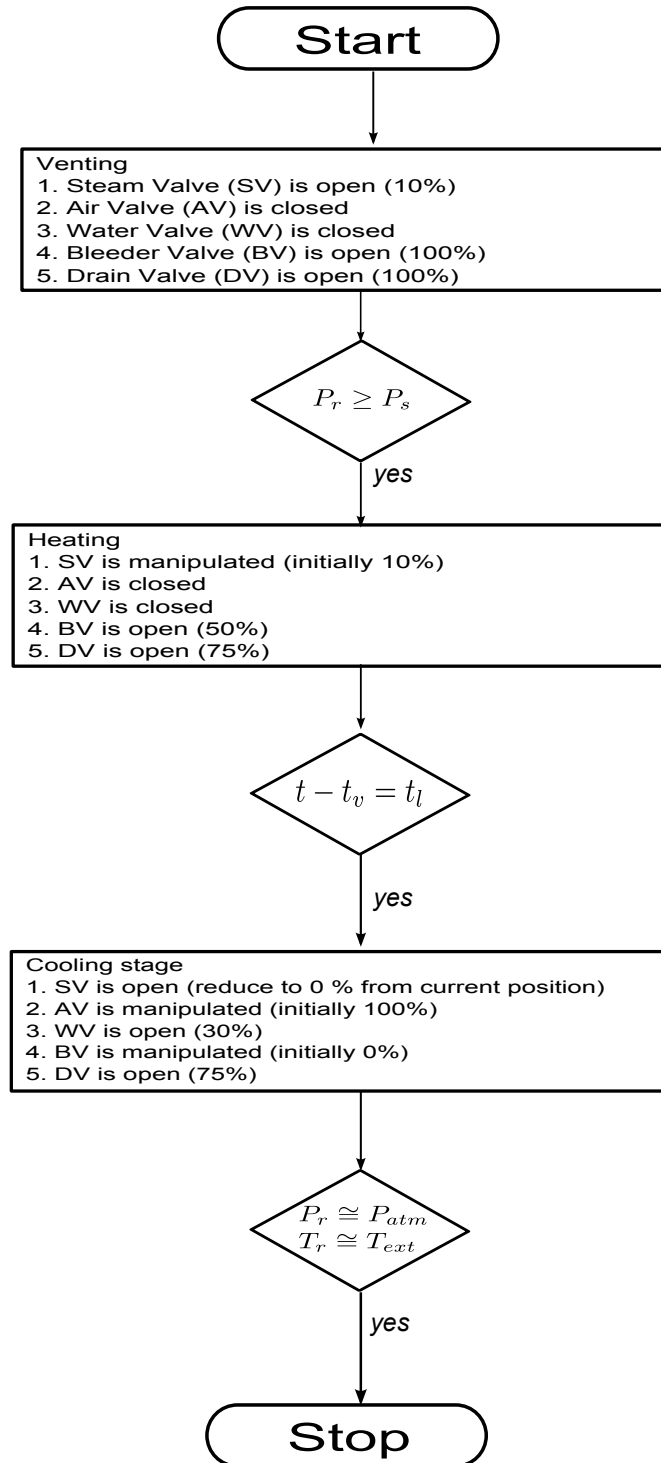
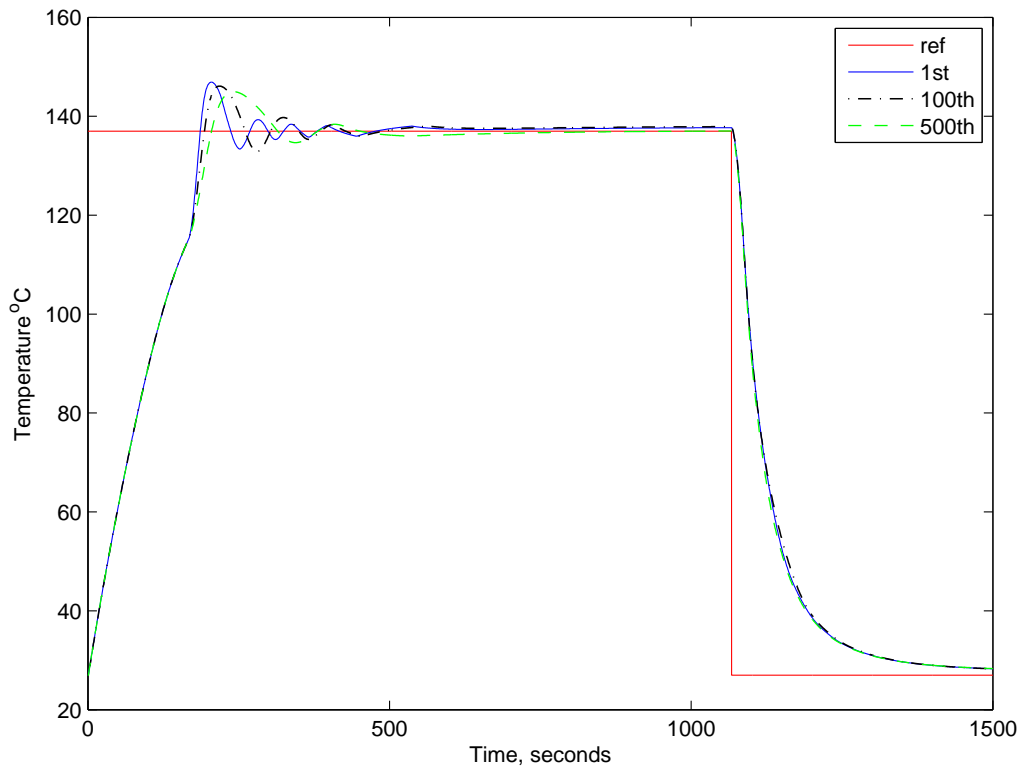
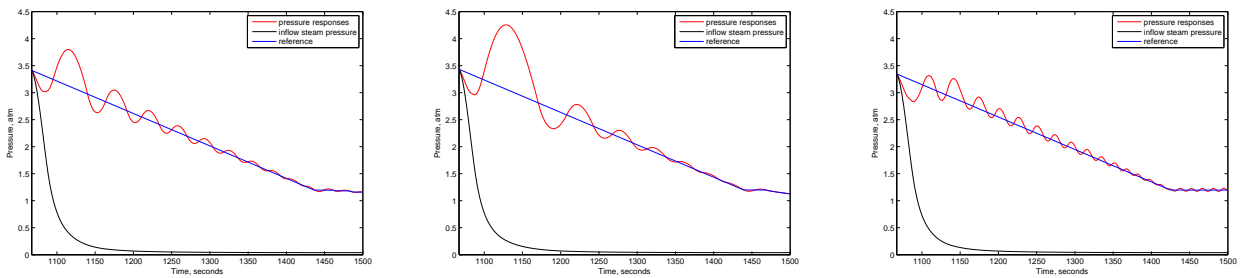


Fig. 6. Control logic implementation:  $P_r$ ,  $P_s$  and  $P_{atm}$  are retort, steam and external pressures,  $t_v$  is the starting time of Heating,  $t_l$  is lethality time,  $T_r$  and  $T_{ext}$  are retort and ambient temperatures.



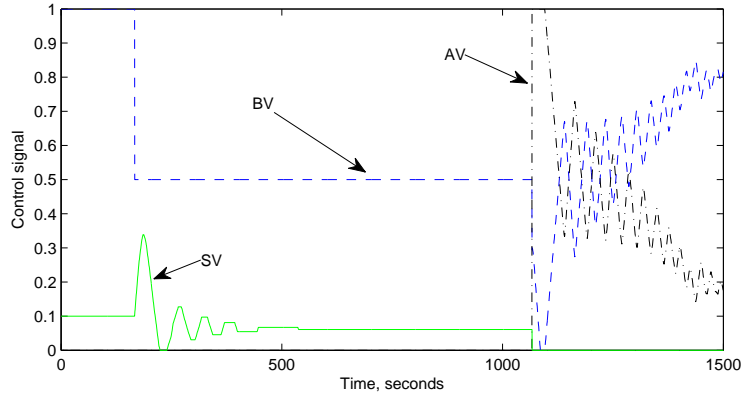
(a) Evolution of the temperature at episodes 1, 100 and 500 (Heating from 200s to 1050s; Cooling from 1050s)



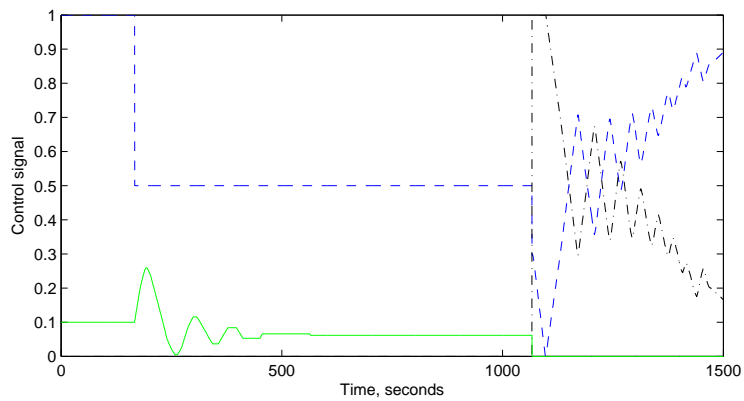
(b) Detail of evolution of pressure during Cooling at episodes 1 (left), 100 (center) and 500 (right)

Fig. 7. Evolution of temperature and pressure during learning on the virtual plant

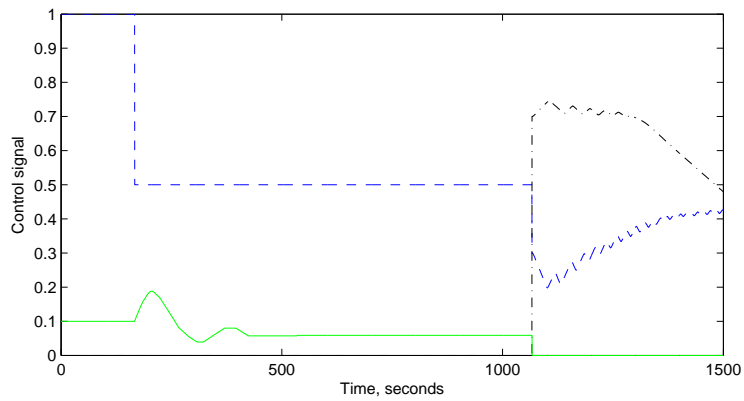




(a) 1st episode



(b) 100th episode



(c) 500th episode

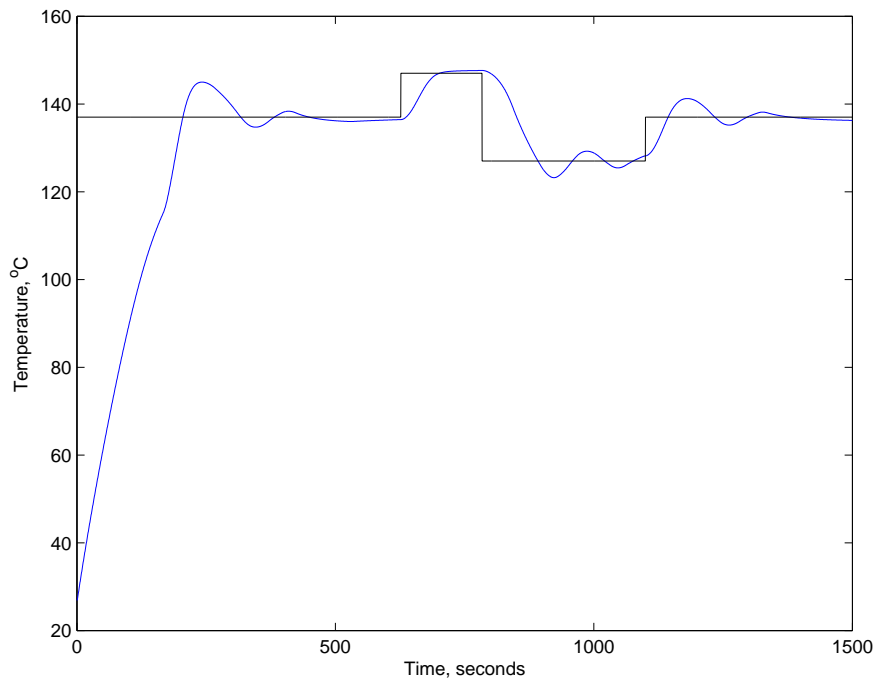
Fig. 8. Control signals during learning on the virtual plant (Heating from 200s to 1050s; Cooling from 1050s)

Table 1  
Heating Control Parameters

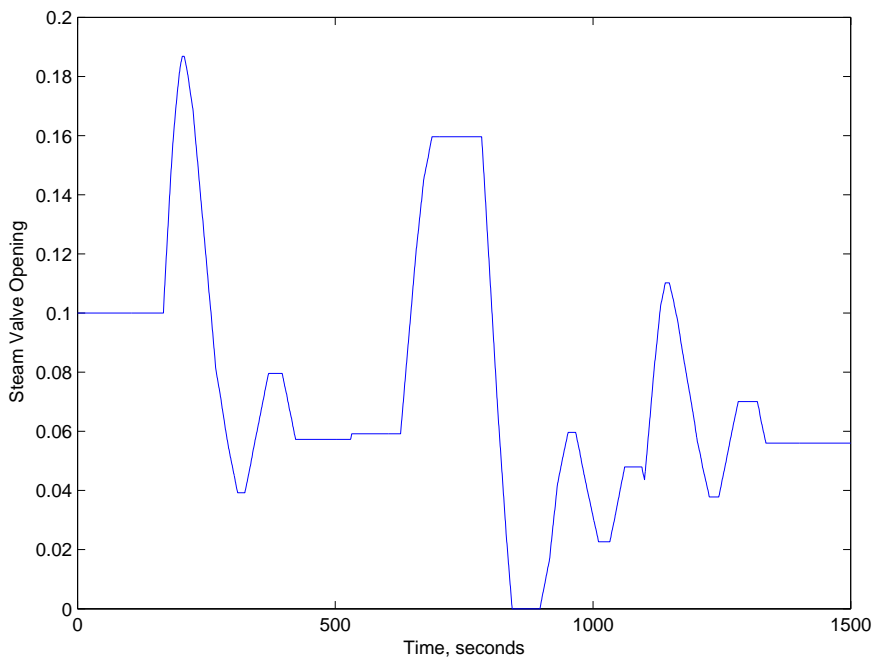
parameters	value	units
learning rate, $\alpha$	0.1	-
forgetting factor, $\gamma$	0.98	-
number of states, $2h + 1$	21	-
span of goal state, $d$	2	$^{\circ}\text{C}$
limited error exploration, $h$	29	$^{\circ}\text{C}$
overlapping degree, $\beta$	5	-
wait action, $a_w$	801	-
controller gain, $k$	$1 \times 10^{-5}$	-
upper limit, $\overline{\Delta u}$	0.008	kg/s
lower limit, $\underline{\Delta u}$	0.0001	kg/s

Table 2  
Cooling Control Parameters

parameters	value	units
learning rate, $\alpha$	0.1	-
forgetting factor, $\gamma$	0.98	-
number of state, $2h + 1$	21	-
span of goal state, $d$	100	Pa
limited error exploration, $h$	$1 \times 10^4$	Pa
overlapping degree, $\beta$	10	-
wait action, $a_w$	601	-
controller gain, $k$	$1 \times 10^{-5}$	-
upper limit, $\overline{\Delta u}$	0.006	kg/s
lower limit, $\underline{\Delta u}$	0.0001	kg/s

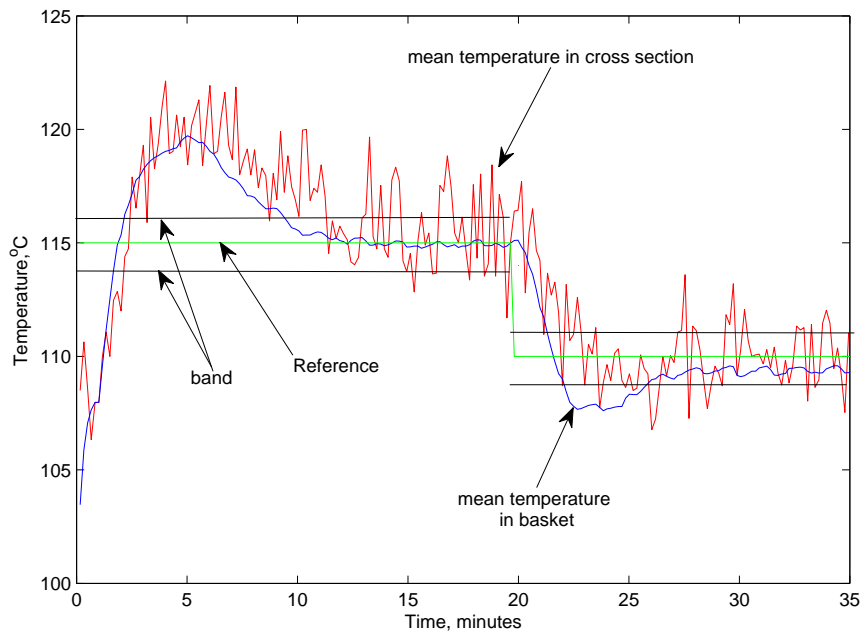


(a) Temperature

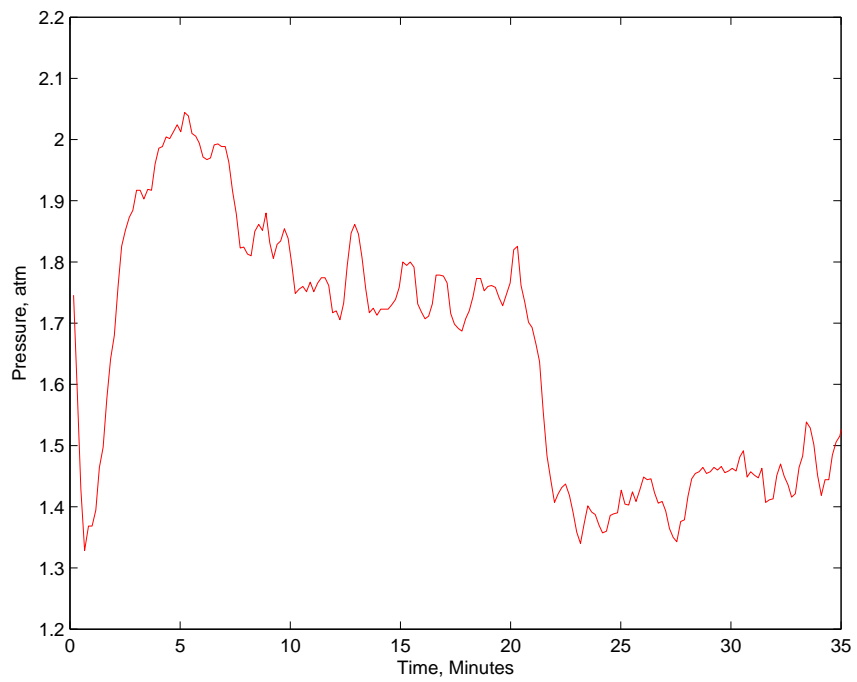


(b) Control signal

Fig. 9. Temperature responses under changes in the temperature setpoint during Heating



(a) Temperature



(b) Pressure

Fig. 10. Temperature and pressure measured in the laboratory plant during Heating, using the proposed control strategy

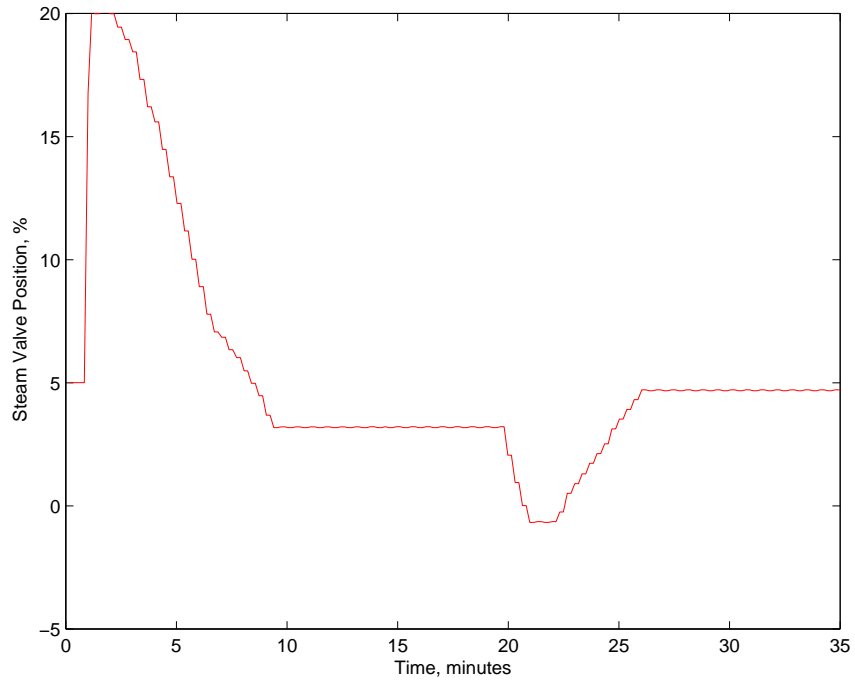


Fig. 11. Steam valve signal calculated by the controller for the experiment in Fig. 10

Table 3  
Control Performances

index	parameters
Time-to-target	2 minutes
Settling time	7 minutes
Maximum overshoot	3 °C